

(19) 日本国特許庁 (J P)

(12) 公表特許公報 (A)

(11) 特許出願公表番号
特表2001-518665
(P2001-518665A)

(43) 公表日 平成13年10月16日 (2001.10.16)

| (51) Int.Cl. ⁷ | 識別記号 | F I | テーマコード* (参考) |
|---------------------------|-------|----------------|--------------|
| G 0 6 F 15/177 | 6 7 8 | G 0 6 F 15/177 | 6 7 8 C |
| 11/20 | 3 1 0 | 11/20 | 3 1 0 A |
| 13/00 | 3 0 1 | 13/00 | 3 0 1 P |
| 15/173 | | 15/173 | Z |
| H 0 4 L 1/22 | | H 0 4 L 1/22 | |

審査請求 有 予備審査請求 有 (全 58 頁) 最終頁に続く

(21) 出願番号 特願2000-514214(P2000-514214)
(86) (22) 出願日 平成10年10月1日 (1998.10.1)
(85) 翻訳文提出日 平成12年3月31日 (2000.3.31)
(86) 国際出願番号 PCT/US98/20532
(87) 国際公開番号 WO99/17217
(87) 国際公開日 平成11年4月8日 (1999.4.8)
(31) 優先権主張番号 08/943, 049
(32) 優先日 平成9年10月1日 (1997.10.1)
(33) 優先権主張国 米国 (US)

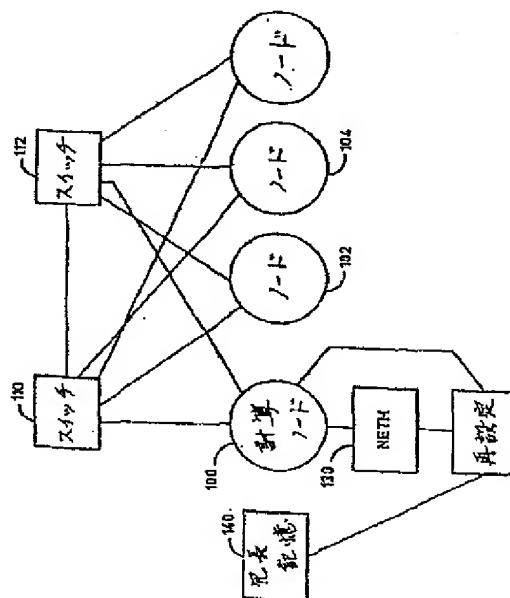
(71) 出願人 カリフォルニア・インスティテュート・オブ・テクノロジー
CALIFORNIA INSTITUTE OF TECHNOLOGY
アメリカ合衆国91125カリフォルニア州パサディナ、イースト・カリフォルニア・ブールバード1200番
(72) 発明者 ブルック, ジョシュア
アメリカ合衆国, 91011 カリフォルニア州, ラ・カナダ, プランブルウッド・ロード, 5657
(74) 代理人 弁理士 深見 久郎 (外5名)

最終頁に続く

(54) 【発明の名称】 分散型計算ノードの高信頼アレイ

(57) 【要約】

頑健な分散型サーバシステムを提供するために冗長記憶および冗長通信を利用するシステム。



【特許請求の範囲】

【請求項1】 冗長分散型ネットワークシステムであって、

複数のシステムノードを含み、前記システムノードの各々は少なくとも2つの通信装置および記憶装置を有し、前記記憶装置はネットワークのための情報の冗長記憶を有し、さらに

前記システムノードの前記通信装置に、いずれか1つのシステムノード内の前記通信装置の各々が複数のスイッチング装置の異なる1つに接続されて、前記システムノードの各々が少なくとも2つの異なる経路のうちの1つを介して互いに通信することができ、したがって冗長通信を提供するように接続される、複数のスイッチング装置と、

ネットワークシステムの動作を妨害し得るシステム機能状態を検出する、検出ルーチンと、

ネットワークの冗長性を利用してシステム機能状態を補償するよう、前記通信冗長性または前記記憶冗長性のうち少なくとも1つを使用してネットワークを再設定する、論理ネットワークプロセスとを含む、システム。

【請求項2】 前記検出ルーチンは、前記システムノードの少なくとも複数上のランを検出し、前記システムノードの各々上では各前記検出ルーチンによって同じプロトコルが実行される、請求項1に記載のシステム。

【請求項3】 前記システム機能状態は、ネットワーク通信内の障害、メモリ記憶内の障害、または、望ましくない結果をもたらす他の種類の障害を含む、請求項1に記載のシステム。

【請求項4】 前記システム機能状態は、ネットワーク通信における障害を含み、前記論理ネットワークプロセスは、接続が異なる接続に変更されるよう指令する、請求項3に記載のシステム。

【請求項5】 前記システム機能状態がメモリ記憶における障害を含むときには、前記論理ネットワークプロセスは、所望の情報が前記冗長データ記憶から得られるよう指令する、請求項4に記載のシステム。

【請求項6】 前記検出ルーチンは各ノードで動作して、ネットワーク内の他のノードに対する接続状態を見る、請求項1に記載のシステム。

【請求項7】 前記検出ルーチンは、ネットワークの状態に関するヒントを使用して前記接続状態を判断するよう動作する、請求項6に記載のシステム。

【請求項8】 前記ヒントは、所定の間隔で各前記ノードによって生成されるハートビート信号を含み、前記検出ルーチンは、前記ハートビート信号を受信して前記ヒントの1つとしての前記ハートビート信号の存在または不在を検出するよう動作する、請求項7に記載のシステム。

【請求項9】 トークンパッシングシステムをさらに含み、各ノードは監視されるチャンネルを介して監視されるノードにおけるイベントを判断し、前記イベントを示すよう前記監視されるチャンネルを介して前記監視されるノードにトークンを渡し、前記監視されるノードは、前記イベントに基づいた動作を示すよう前記トークンを戻し、各ノードは、該ノードとは異なるノード上で対応するイベントが起きずとも該ノードにおいて起こすことのできるイベントの数を制限するよう、所定数のトークンのみを有する、請求項7に記載のシステム。

【請求項10】 各前記ノードがネットワークの同じ履歴を見ることを保証する手段をさらに含む、請求項6に記載のシステム。

【請求項11】 前記接続は、計算ノードのグループが分離されることのないように行なわれる、請求項1に記載のシステム。

【請求項12】 前記スイッチは、可能な限り最も非ローカルな方法で前記ノードを接続する、請求項1に記載のシステム。

【請求項13】 前記スイッチは、互いに最も遠い2つのノード間を接続する、請求項12に記載のシステム。

【請求項14】 前記接続は、どの2つのノードの故障によってもノードの1グループをノードの他のグループとの通信から分離することのないように行なわれる、請求項13に記載のシステム。

【請求項15】 各ノードは、少なくとも2つの経路によって各他ノードと接続され、

ネットワークモニタをさらに含み、該ネットワークモニタは、各前記ノードにおいて稼動して、前記ネットワークモニタがその上で稼動するローカルノードから各リモートノードへの各接続経路上のすべての接続を監視する、請求項1に記

載のシステム。

【請求項16】 高信頼ユーザデータプロトコルをさらに含み、該プロトコルは、前記ローカルノード上で稼働し、該ローカルノードから他のノードへの通信リクエストを受信し、かつ、前記ネットワーク監視プロセスから経路を判断する、請求項15に記載のシステム。

【請求項17】 ノード間の物理的な接続を、異なるノード接続に変更することを可能にする、論理ネットワーク相互接続を使用した前記通信経路の再設定をさらに含む、請求項16に記載のシステム。

【請求項18】 前記ノード間の動作可能な接続を判断するネットワークモニタと、稼働中のノードのための情報を処理する高信頼ユーザプロトコルと、前記動作可能な接続に基づいて通信を再設定する論理ネットワークと、をさらに含む、請求項1に記載のシステム。

【請求項19】 前記記憶装置は、各ノードの各ディスク上の情報の一部のみを記憶する、請求項1に記載のシステム。

【請求項20】 各ノードの各ディスクは、他のディスク上の情報の特性を示す情報もまた記憶する、請求項19に記載のシステム。

【請求項21】 冗長分散型サーバであって、
分散型計算ノードのアレイを含み、前記計算ノードの各々は各他ノードとは異なる情報を記憶し、前記記憶される情報は前記計算ノード間で冗長となっており、さらに

前記計算ノードのアレイに接続され、前記計算ノードのアレイ間に冗長通信経路を提供し、どのような所定数のネットワーク障害もシステムの残りのノードの動作に影響を及ぼすことのないように動作する、スイッチングシステムを含み、

前記計算ノードの各々はネットワークステータスを判断する同じプロトコルを実行し、よって、各前記計算ノードが同じネットワーク履歴を見るようにする、サーバ。

【請求項22】 前記各ノード上に記憶される情報は、所望の情報のすべてではない一部のみを記憶し、どの2つのノードも同じ情報を記憶することはない、請求項21に記載のサーバ。

【請求項23】 前記記憶される情報は、情報部分および冗長部分を含み、前記冗長部分は他のノードのための情報部分のみを示す情報である、請求項22に記載のサーバ。

【請求項24】 前記冗長部分は、複数の前記ノードが前記情報部分を形成するようアレイへと配列されているアレイ符号から形成され、かつ、前記冗長部分は、前記アレイの対角線方向に沿ったチェックサムで形成される、請求項23に記載のサーバ。

【請求項25】 集合的にシステムデータを記憶する複数のノードを含み、各ノードは生のデータおよび、前記各ノードとは異なる他のノード内に記憶されている生のデータを示す冗長データを記憶し、さらに

前記複数のノードからの分散型の読出を実行するよう動作する制御プロセスを含み、前記分散型読出は、ノードの可用性に関するパラメータを判断するステップと、もし前記パラメータが可用性を示す場合には前記複数のノードから前記生のデータを読出すステップと、もし前記パラメータが可用性を示さない場合には、前記複数のノードよりも少ない数のノードから前記生のデータおよび前記冗長データの両方を読出すステップとを含む、冗長ネットワーク。

【請求項26】 誤り訂正符号を示す生の情報および冗長情報を複数の情報ノードに記憶するステップと、

前記情報ノードのユーザビリティを示すパラメータを判断するステップと、

前記パラメータが前記複数のノードが使用可能であることを示す場合には前記複数のノードから前記生の情報を読出し、かつ、前記パラメータが前記複数のノードのうち少なくとも1部分が使用可能ではないことを示す場合には前記複数のノードよりも少ない数のノードから前記生のデータおよび前記冗長データの両方を読出すステップとを含む、誤り訂正符号を使用する方法。

【請求項27】 ノードからの情報を表わす、アレイの各列を形成することによって、情報のアレイを形成するステップと、

データを示す生の情報を含む、各列の生の部分を形成するステップと、

冗長性に関する情報を示す、冗長情報を形成するステップとを含み、前記冗長情報は、前記各ノード以外の他のノードに関する情報を示し、該情報は、前記他

のノードから情報を得る特定の形状のエンベロープに沿って取られる、ノードのアレイのための符号化方法。

【請求項28】 前記エンベロープは、前記アレイの縁を超えて他ノードに延びる対角線である、請求項27に記載の符号化方法。

【請求項29】 アレイの列に各ノードをマッピングするステップと、
前記アレイの列から冗長情報の2つの行を形成し、前記2つの行を前記列内に位置付けることによって、 $N-2 \times N$ 個の情報シンボルおよび $2 \times N$ 個の冗長情報シンボルを含む、結果として得られる $N \times N$ のアレイを形成するステップとを含み、前記パリティシンボルは以下の式に従って構築される、複数のノードで形成される冗長符号化ネットワークを形成する方法。

【数1】

$$C_{n-2,i} = \sum_{K=0}^{n-3} C_{k,(i+k+2)_n}$$

$$C_{n-1,i} = \sum_{K=0}^{n-3} C_{k,(i-k-2)_n}$$

式中、 $i = 0, 1, \dots, n-1$ および $\langle x \rangle_n = x \bmod n$ である。

【請求項30】 冗長ビデオサーバシステムであって、

複数のシステムビデオ記憶ノードを含み、前記システムビデオ記憶ノードの各々は、少なくとも2つの通信装置および記憶装置を含み、前記記憶装置はビデオ情報の冗長記憶を含み、さらに

前記システムビデオ記憶ノードの前記通信装置に、いずれか1つのシステムノードにおける前記通信装置の各々が複数のスイッチング装置の異なる1つに接続されて、前記システムノードの各々が少なくとも2つの異なる経路のうち1つを介して互いに通信することができ、したがって、冗長通信を提供することができるようになる、複数のスイッチング装置と、

特定の数の故障がビデオサーバの動作を中断することのないような形をとる、前記スイッチング装置およびビデオ情報の冗長記憶上の接続と、

ネットワークシステムの動作を妨害し得るシステム機能状態を検出する、検出ルーチンと、

ネットワークの冗長性を使用してシステム機能状態を補償するよう、前記通信冗長性および前記記憶冗長性のうち少なくとも1つを使用してネットワークを再設定する、論理ネットワークプロセスとを含む、システム。

【発明の詳細な説明】**【0001】****【技術分野】**

本願は、頑健な通信ならびに分散型の読出および書込動作を形成するような形で、情報の冗長通信および記憶を含むネットワークを形成する、分散型計算ノードの高信頼アレイに関する。また、該システムは、冗長性が必要であることを示す状態を検出し、そのような状態に対応して該状態に対する補償をするために再設定を行なうこともできる。

【0002】**【背景技術】**

分散型環境上の計算および記憶を使用すれば、既存のハードウェアおよびソフトウェアを増強することのできる可能性は大いにある。このようなシステムは、分散型の、可用性に優れた記憶サーバとして使用することができるであろう。可能な用途としては、マルチメディアサーバ、ウェブサーバ、およびデータベースサーバ等がある。しかし、より一般的には、この種のシステムは、情報を複数の場所で分散する必要がある、どのような応用にも使用することができる。

【0003】

しかしながら、コストを過度に増すことなく高い信頼性を提供するように、接続、監視および動作を正しく組合せることは、非常に困難である。

【0004】

障害を補償することのできる冗長記憶システムを提供する方法はよく知られている。このようなシステムの一例が、「RAID (reliable array of independent disks)」と呼ばれるシステムである。RAIDタイプのシステムの2つの例が、米国特許番号第5, 579, 475号および第5, 412, 661号に記載されている。これらのシステムは冗長データ記憶を提供することで、システムのどのディスクの故障も、該システム内の他の場所にある冗長データによって補償されるようにしている。

【0005】

システム内の各コンピュータ（「ノード」）が他のノードと接続されている通

信システムはよく知られている。その一例がイーサネットであるが、これはバス型プロトコルであって、計算ノードがバスを介して通信する。サーバは典型的に、すべてのノードのための共用データをすべて格納している。これらのノードは、ローカルなデータ記憶もまた有してもよい。

【0006】

単一のネットワークシステムは、ノードとサーバとをつなぐ単一のイーサネットリンクを含む。したがって、サーバへの接続もしくは通信に、またはサーバ自身の内部になんらかの障害が発生すると、ノードはもはやサーバから従来のデータアクセスサービスを得ることができなくなる。それらノードは、その後はスタンドアロンモードで動作することを強いられ、したがってローカルに利用可能なデータを使用して動作することしかできなくなる。

【0007】

このようなシステムの信頼性を高めることを目的としたサーバ型システムが知られている。このようなシステムの1つは、デュアルバス接続を使用する。各計算ノードには、2つの別個のイーサネットカードを使用して2つの別個のバスを介して2つの別個のサーバにつながる、2つのイーサネット接続が提供される。これは、実質的には、各々がハードウェアおよび記憶の完全な予備を有する、2つの別個のシステムである。

【0008】

この場合、いずれかの接続またはバスにエラーがあっても、他方のバスを介して正常な動作を続けることができる。2つの冗長バスおよび2つの冗長サーバを有するシステムを、デュアルバス、デュアルサーバと呼ぶ。このようなデュアルバス、デュアルサーバシステムは、単一のネットワーク障害は許容することができる。しかし、このようなシステムでは、通常、各サーバ上にすべての情報の複製を作成しなければならない。

【0009】

【発明の開示】

本願に記載されるシステムは、パーソナルコンピュータ等の比較的ローパワーのワークステーションを使用して、既存のハードウェアおよびソフトウェアを増

強するものである。これらのパーソナルコンピュータは、冗長接続によって接続される。この接続は、既存のハードウェア、たとえばローカルおよび／またはワイドエリアネットワークを利用することができる。

【0010】

本願は、分散型計算ノードのアレイで形成された、冗長分散型サーバを開示する。それら計算ノードの各々は、特別な冗長方式で情報を記憶し、また、頑健な通信を保証するプロトコルを実行する。

【0011】

該システムは、好ましくは、ある特定の数のネットワーク障害が、システムの残りのノードの動作に影響を及ぼすことのないように、ネットワークにフォールトトレランスをもたらす、特別なアーキテクチャおよび動作を含む。しかし、そのノードはどれ一つをとっても、すべての情報を二重に記憶することはない。

【0012】

該サーバシステムは、冗長通信および冗長記憶を含む。冗長通信は、各ノードが少なくとも2つの異なる経路のうちの1つを介して他の各ノードと通信することができるようにする、システムアーキテクチャによって得られる。冗長記憶は、特別な冗長符号化方式を使用して情報を冗長に記憶することによって得られる。

【0013】

該サーバシステムはまた、システムの機能状態を検出する分散型検出ルーチンを実行する。システムの機能状態としては、たとえば、ネットワーク障害がある。ネットワーク障害は、リンク断、または動作不能ノードもしくは動作不能スイッチングデバイス等の、通信障害を含み得る。しかし、より一般的には、システム機能状態は、ネットワークの動作を妨げるおそれのある、どのような状態も含み得る。システム機能状態は、システムの冗長性によって補償することができる。

【0014】

該サーバシステムは、好ましくは、システムの機能状態を検出するネットワーク監視プロセスを実行する。論理ネットワークプロセスがシステムを再設定して

、冗長性を利用してシステム機能状態を補償するようにしている。

【0015】

該システムはまた、システム障害があった場合に代替的な動作ができるようにする、分散型読出および書込システムを使用する。この代替的な動作は、システムの冗長性を使用する。

【0016】

【発明を実施するための最良の形態】

本発明の目的、利点および特徴は、添付の図面を参照して以下の詳細な説明を読むことにより、より容易に理解されるであろう。

【0017】

図1は、高信頼冗長分散型ネットワークサーバシステムの、最も基本的な第1の実施例を示す。このシステムは、複数の計算ノード（「ノード」）と、それらノードの間でスイッチングを行なうネットワークとで形成されている。

【0018】

図1のネットワークは、ノードとネットワークとの間に、通信および記憶の双方について冗長性を有する。この冗長性は、所定の数のシステム機能状態を補償するのに使用することができる。この冗長性によって補償されるシステム機能状態は、ネットワークにおける障害（「通信障害」）、ディスク、揮発性メモリ、もしくはデータを記憶する他のいかなる種類のメモリであってもよい、メモリの記憶における障害（「メモリ障害」）、または、望ましくない結果をもたらす他のどのような種類の障害をも含み得る。

【0019】

この分散型サーバシステムはまた、検出プロセスを含む。検出プロセスは、各ノードにおいて稼動して、ネットワーク内の他のノードへの接続状態を見る。各ノードは、ネットワークの状態に関するヒントの集積を使用して、同じプロトコルに従ってネットワークを見る。この検出プロセスは、両サイドがネットワークの同じ履歴を見ることができるようになっている。この検出プロセスは、分散されてはいるが、トークンパッシングシステムを使用して、所望のしきい値内で整合する、ネットワークのノード間のネットワーク履歴を維持する。トークンにより

、一方側が、他方側がアクションをとったことを確認せずにとることのできるアクションの数を、所定数に制限することによって、2つの側の自由度が制限される。

【0020】

この検出プロセスは、他のプログラムおよびユーザアプリケーションから見えないところで実行される。検出プロセスの好ましい形態は、監視対象となるシステムに関する情報を集めるよう動作する、ネットワーク監視（「NETM」）プロセスを使用するものである。このNETMプロセスは、好ましくは、他方ノードが正しく動作しているかどうかを判断するが、より一般的には、NETMプロセスは、ユーザビリティに関するパラメータを判断する。このパラメータは、以下に説明するように、システムが動作中かダウンしているかを示す表示を含み得る。このパラメータは、システムがどの程度ビジーであるかを示す表示もまた含み得るが、この表示は、負荷の最適配分に使用することができる。

【0021】

図1のシステムは、2つのスイッチ110および112によって接続される4つの計算ノード（「ノード」）100、102、104および106を用いて、本発明の特徴を示している。各ノードは、2つの異なる、したがって冗長な経路を介して、他の各ノードと通信することができる。たとえば、ノード100は、ノード100とスイッチ110との間の相互接続120を介して、ノード106と通信することができる。また、スイッチ110からノード106への経路122を介して冗長相互接続を与える、完璧に別個の経路が存在する。これとは別に、ノード100は、ノード100からスイッチ112への相互接続124、および、スイッチ112からノード106への相互接続126を使用して、ノード106と通信することもできる。したがって、各ノードは、少なくとも2つの完全に別個の、冗長な接続経路によって、他の各ノードと接続されている。

【0022】

この冗長通信能力は、ある通信リンクの使用を避ける方が好ましい場合に、異なる経路を選択することができるようする。たとえば、スイッチ110、または、接続線120および／または122のどこかの部分が失われた場合にも、スイ

ッチ112を介して、接続線124および126を通じて、通信を行なうことができる。

【0023】

情報もまた、冗長な方法で記憶されており、ネットワークのどこかの部分が故障するか、他の形、たとえば膨大なトラフィック等により利用不可能となった場合にも、どのような情報も検索することができる。この冗長記憶機構を、図1においては要素140として示す。システム内のノードの全個数を n で表わし、選択数を k で表わすとする、冗長記憶140内のデータは、 $n-k$ 個のノードが失われた場合にも、システムから所望のデータを得ることのできる能力に影響が及ぶことのないような態様で、記憶されるのが好ましい。このような記憶は、好ましくは、ノードの各々に冗長情報が記憶される最大距離分離（「MDS」）符号方式に従ってデータを記憶することによって行なわれる。この冗長情報を他のノードのデータとともに使用して、欠けてしまった1つ以上のノードのためのデータを再構築することができる。

【0024】

検出プロセスによって、動作不能ノードや破断した通信リンク等の望ましくないシステム機能状態が検出されると、再設定プロセス140が実行される。この再設定プロセス140は、記憶冗長性または通信冗長性の少なくとも一方を使用することができるというその能力のおかげで、障害に対して頑健である。再設定プロセスは、特定の障害があった場合にもシステムが動作することができるようにする。しかし、このプロセスは、専用のスイッチングを何ら必要としないであろう。たとえば、ノード100と106との間の経路は、120/110/122を介する経路1上で、または、124/112/126を介する経路2上で、構築することができる。正常動作中には、通信は、経路1、次に経路2、その次に経路1、等のように、交互に行なわれる。しかし、経路1内に障害または過負荷があれば、すべての通信は経路2上で行なわれるようになる。これは、通信が適切に方向付けられるという意味で、再設定と言える。いずれにせよ通信の半分が経路2を介するように方向付けられていたであろうとはいうものの、この再設定によって、すべての通信が経路2を介するように仕向けられるのである。

【0025】

したがって、図1は、本願において説明する分散型サーバの基本的な特徴を示している。これらの特徴は、通信の冗長性、記憶の冗長性、その冗長性によって補償することのできるイベントの検出、および、そのイベントを補償するために冗長性を使用する再設定、を含む。

【0026】

冗長通信

図1のシステムは、4つのノード100～106ならびに2つのスイッチ110および112を含む、簡単な冗長接続を示す。これらのノードは、好ましくは、各々が2つのPCIバス型通信カードを有する、パーソナルコンピュータ（「PCS」）等のスタンドアロンワークステーションである。通信カードは、スイッチを介して、他のPCSの同様の通信カードと通信する。通信カードのプロトコルは、イーサネットまたは他の市販のどのような種類のものであってもよい。好ましいシステムは、図2に示すスイッチングノード200に対して、ミリネット（Myrinet）スイッチを使用する。ミリネットスイッチは市販されており、また、ボーデン（Boden）他による「ミリネット：ギガビット／秒のローカルエリアネットワーク（"Myrinet: a gigabit per second local area network"）」IEEE Micro 1995、に記載されている。

【0027】

本発明によって使用される特別なノード接続は、ネットワーク通信障害があった場合にも正常に動作する能力を改善する、通信冗長性を提供する。ネットワーク通信障害は、スイッチ障害、リンク断、またはスイッチ故障等の、障害を起こした通信を含む。接続は、単一の通信障害または通信障害の組合せが通信の破壊またはノードの孤立化を引起す可能性を最小に抑えるような方法で構築される。正しい接続の重要性を、以下に説明する。

【0028】

図2は、4つのスイッチ220～226を使用して8つの計算ノード200～214を接続するシステムを示す。各計算ノードは、使用可能な2つの相互接続リンク経路を含み、これにより、通信の冗長性が与えられている。

【0029】

しかし、図2のシステムにおける通信故障は、計算ノードのグループを意に反して「分離する」おそれがある。これらの孤立した計算ノードのグループは、それらが分散型サーバの他の稼動中のノードのすべてともはや通信することができないという意味で、分離されている。

【0030】

1例として、もしスイッチ224および226が両方とも故障してしまえば、計算ノード200～206が計算ノード208～214とは完全に分離してしまう。この場合、使用可能ではあるがあまり好ましくない状態の、分離し得るシステムができることになる。

【0031】

たとえば、使用されるMDS符号が、データを再構築するのに6個から8個のノードを必要とするものとする。もしこのシステムが上述のように分離されれば、半数のノードのみしか通信ができなくなる。通信可能なノードは4つなので、この場合は障害によって、データの再構築ができなくなってしまう。

【0032】

図3のような接続構造が好まれる。この10ノード、4スイッチを含むシステムでは、通信障害がある場合の相互接続が改善されている。コネクションインターフェイスは、いずれか2つのスイッチが失われても、最悪の場合でも2つの計算ノードのみにしか影響が及ぶことのないように作られている。たとえば、スイッチ320および326が故障した状態を示す図4を参照されたい。太い線は、この故障によって影響を受ける通信線を示す。この2つのスイッチの故障によって分離されるのは、計算ノード304および312のみである。他のすべてのノードは、全く問題なく動作でき、どれも分離されることはない。

【0033】

フォールトトレランスの重要な部分は、スイッチおよびノードの特定の相互接続によって得られる。上述の例のように、図2のシステムは、計算ノードの半分ずつが分離されてしまうおそれがある、という欠点を有する。分離されたシステムは計算ノード200～206を含み、それらのノードは通信することはでき

るが、ノード208～214のグループとは分離されている。

【0034】

この問題の別の例を、図12に示す。図12は、複数のスイッチングノードを使用して、多数の計算ノードを相互接続する方法の1例を示す。各スイッチングノードNは、隣接する2つの計算ノードCの間に位置付けられる。これは、使用可能ではあるが、あまり好ましくない構成である。計算ノード1200および1202が同時に障害を起こしてしまった場合、このシステムの通信能力は、図12に示した破線に沿って分割されてしまうことに留意されたい。これにより、システムの半分1204が、システムの他方の半分1206とは実質的に分離されてしまう。

【0035】

本願に開示した接続の目的は、2つの通信故障によって形成されるおそれのある、この種の分離を防ぐことである。好ましいシステムは、可能な限り最も非ローカルな方法でノードを接続する方法を説明する。これは、各スイッチングノードが最も近い2つの計算ノードに接続される図12のシステムとは、対称的である。発明者は、非ローカルなスイッチ間を接続するこの非自明のシステムが、非常に優れたフォールトトレランスをもたらすことを発見した。

【0036】

図13はそのような装置を示す。各ノードは2つのスイッチに接続されており、この図は、最も遠い2つのスイッチ間が接続されている様子を示す。図13のようにスイッチおよびノードを並べて図示した場合、互いに対して物理的に最も遠い2つのスイッチ間を結ぶ接続は、直径として表わされる。このような接続は、どれか3つのスイッチが失われた場合にも、ユニット全体の半分ずつを分離させるような影響が生ずることはない、という利点を有する。逆に、いずれか2つの場所でユニットに破損が生じた場合にも、多くのノードの間で通信は依然として行なうことができる。どこか3つの場所で損失があった場合でも、システム内のノードの全個数に関係なく、ある一定数のノードー直接影響を受けるノードーのみが分離されるにすぎない。

【0037】

たとえば、ロケーション1310で通信故障が生じ、ロケーション1312で別の破損が生じたとする。スイッチ1300がスイッチ1304を介して依然としてスイッチ1302に接続しているので、ノードが依然として通信することができるのは明らかである。スイッチ1300はまた、スイッチ1308を介してスイッチ1306にも接続されている。同様に、これらのスイッチはすべて、破損があった場合にも互いに接続されている。さらに、この好ましいシステムにおいては、必要とされるであろうノード間接続は、最大でも、システム全周の1/4である。

【0038】

この非ローカルという概念は、リング型以外の構成にもあてはまる。たとえば、リングと見なすことのできる構成であれば、どのようなものも代替的に使用することができる。

【0039】

図1から図3に示した好ましいサーバシステムは、ミリネット相互接続技術を使用して冗長ネットワークを介して接続される、パーソナルコンピュータ型のワークステーションを使用する。もちろん、代替的に、100MBイーサネット等の他の通信技術を使用することもできる。これらのシステムはすべて、共通に、障害を起こしたリンクがあった場合にも冗長性を維持するという能力を有する。このシステムは、どのような数の通信要素とも、ともに使用することができる。ただし、好ましい数、また開示される数は、2である。

【0040】

冗長記憶

図1の好ましい実施例において、各ノードは、所与の記憶データのうち、1部分のみを記憶している。この記憶データは、ローカルノード内に実際に記憶されている各情報の一部、および他のノードが記憶している一部を使用して、取出される。この概念を図5を参照して説明する。図5は、ビデオサーバを示す。開示されるサーバは、図示されるように表示されるビデオを示すデータを提供する。各計算ノードは、図示するように、データ全体の半分の量を記憶している。このデータは、そのデータを要求する1ノード内に記憶されているデータを、他のノ

ード内に記憶されているデータと組合せることで、どのビデオフレームも再構築することができるような方法で、冗長に記憶されている。

【0041】

この記憶方式によれば、どのノードも、そのノードが他のすべてのノードから分離されていない限り、そのノードが所望する情報を受取ることができる。この方式は、分散型サーバ内に多数の故障が生じた場合に備える、記憶冗長性を提供する。

【0042】

しかし、より一般的には、ここに規定する好ましい方式は、全部で n 個のノードのうち、 k 個の稼働中のノードのサブセットから、データを再構築することができるようにするものである。 $k=2$ および $n=4$ の場合の例を、以下に示す。

【0043】

図6は、計算ノードの1つに故障が生じた場合に、送り出されたビデオのどの項目でも、残りの計算ノードで復元することができる方法を示す。これは、冗長記憶を可能にするどのような符号方式でも達成することができる。

【0044】

好ましいシステムは、いずれか2つの通信リンクを失っても、そのサーバシステムの他のいかなる通信機能も失うことがなく、かつ、障害を実際に含むノード以外の他ノードに影響を及ぼすことがない。

【0045】

このシステムの冗長記憶特徴は、各ノード内に、データ全体よりも小さいサイズ、全データの半分の量の、符号化されたデータを記憶する。したがって、 k 個の稼働中のノードを有するシステムにおけるサイズ K の各ファイルについて、この好ましい実施例においては、そのファイルの K/k が、サーバの各ノード上に記憶されている。ファイルの他の $(k-1)$ は、他の $k-1$ 個の稼働中のノードから得られる。

【0046】

X符号

記憶の冗長性は、好ましい実施例に従えば、ノード間で情報の記憶を分散させ

ることによって得られる。上述のように、サイズ K の情報の各項目について、好ましいシステムは、各ノード内に K/k のデータ（情報のオリジナルサイズ）を記憶する。ここで、 k は、データを再構築するのに必要となるであろうノードの数を表わす。各ノードは、他のいずれかのノードに記憶される情報の他の K/k にアクセスすることによって、情報のいずれの項目も再構築することができる。情報は、好ましくは、その情報を記憶するために最大距離分離（「MDS」）符号を使用して記憶される。情報を記憶するための好ましい形態は、 X 符号と呼ばれる新しい符号化システムを使用する。ここに記載する X 符号は、ノード間、より特定的にはノードのディスク間に分散される、情報の各項目を記憶するための、特別な、最適化された符号である。

【0047】

最も好ましくは、情報の一部のみ、符号化されたデータのある部分のみが、各ノード上に記憶される。各ノードは、他のノード上の情報のなんらかの特性を示す情報もまた記憶する。この特性とは、たとえば、他のノード上のデータの合計を示す、チェックサムまたはパリティであり得る。この情報は、他のノード上の情報を再構築するために、それら他のノード上の情報とともに使用される。

【0048】

上述のように、使用される好ましい符号は X 符号であって、これを以下に詳細に説明する。 X 符号は、 $N \times N$ の最大距離分離（「MDS」）アレイ符号であって、 N は好ましくは素数である。この符号は、排他的論理和（「XOR」）および循環シフト演算のみを使用して、符号化および復号化することができる。このため、 X 符号は、リードソロモン符号等の計算がより複雑な符号に比べて、はるかに高速で符号化および復号化することができる。

【0049】

この X 符号は、最小列距離3を有する。これは、この符号が1つの列誤りまたは2つの列消失を訂正することができることを意味する。 X 符号は、 X 符号内の単一の情報ビットまたはシンボル等の、単一の情報ユニットの変化が、常に、2つのパリティビットまたはシンボルだけにしか影響を及ぼさない、という、特定の特性を有する。したがって、情報が更新されるたびに、変更が必要となるの

は、それら2つのパリティビットまたはシンボルのみである。

【0050】

X符号のシステムは、図15に示すアレイを使用する。各列1500は、単一ノード内の情報を表わし、各ノードにマップされる。パリティシンボルは、列ではなく行に記憶される。

【0051】

この符号は、ネットワークのすべてのノードを使用して配列されて、集合的に $N \times n$ のアレイを形成する。ここで、好ましくは $N = n$ である。このアレイは、 $N-2 \times N$ の情報シンボル、および $2 \times n$ のパリティシンボルを含む。図14 (A) は、 $n = 5$ の場合のアレイの例を示す。ノード1400の部分は情報を表わし、四角で囲った要素は、1ビット、1セクタまたはディスクの他の何らかのユニット等の情報の1ユニットを表わす。本明細書では、これらのユニットを総称的にシンボルと呼ぶ。

【0052】

非情報部分1402は、冗長情報を表わす。ここで説明するように、アレイの単一系列によって表わされるディスク、たとえばディスク番号1404のディスクに対して、冗長情報1402は、他のディスクからの冗長情報を表わす。すなわち、この冗長情報は、1404以外のディスクからのみ得られるものである。

【0053】

このX符号システムは、ディスク1404全体の内容を表わす列を形成する。X符号のパリティシンボルは、ディスク上の2つの付加的な行1402から形成されている。したがって、各ディスクは、 $N-2$ 個の情報シンボル、および2個のパリティシンボルを有する。1列内の1シンボルの誤りまたは消失は、列消失から復元することができる。

【0054】

ここで具体的な符号化手順を説明する。 i 番目の行および j 番目の列のシンボルを C_{ij} とすると、X符号のパリティシンボルは、以下の式に従って構築される。式1:

【0055】

【数2】

$$C_{n-2,i} = \sum_{K=0}^{n-3} C_{k,(i+k+2)_n}$$

$$C_{n-1,i} = \sum_{K=0}^{n-3} C_{k,(i-k-2)_n}$$

式中、 $i = 0, 1, \dots, n-1$ および $(x)_n = X \bmod n$ である。

【0056】

これは、幾何学的には、傾き1および-1の対角線方向にそれぞれ沿ったチェックサムを表わす、パリティ行と解釈することができる。

【0057】

図14(A)は、第2のパリティチェック行1412が存在しないかまたはすべてゼロ(0)であると仮定して、第1のパリティチェック行1410を得る方法を示す。なお、存在しないかまたはすべてゼロのパリティチェック行は、仮想ゼロ行と呼ぶ。チェックサムは、傾き-1のすべての対角線方向上で形成される。図14(A)においては、三角形をすべて加算することで、第1のパリティチェック行1410が形成される。具体的には、要素1414、1416、1418および1420を加算して、パリティ要素1422を形成するのである。

【0058】

図14(B)は、例示の複数の単一ビットについて、第1のパリティチェック行を計算する1例を示す。対角線方向の要素1414、1416、1418および1420については、 $1+1+1+0$ を計算することでパリティ1が導出され、これがシンボル1422として記憶される。

【0059】

対角線は、アレイの外縁に達すると隣接する行に続く。たとえば、要素1432、1434、1436および1438を含む対角線方向の行1430は、次の行の最上部に続いて、1440で始まる。パリティシンボル1436は、シンボル1432、1434、1438および1440を加算したものに対応する。図14(B)は、これらのシンボルが $0+0+0+1$ に対応して、1となることを示しており、シンボル1436として値1が記憶される。

【0060】

第2のパリティチェック行は、傾き+1の対角線方向から形成される。図14 (C) は、図14 (D) とともに、第2のパリティ行の同様の計算方法の具体例を示す。行1440は、シンボル1442、1444、1446、1448および1450を含む。パリティシンボル1450は、 $1442 + 1444 + 1448 + 1446$ で計算される。図14 (D) は、具体的な例を示しており、ここでは、 $+0 + 0 + 1 = \underline{1}$ から、パリティ0が得られる。

【0061】

図14 (E) は、これら2つのパリティチェック行を組合せることによって形成される、完全な符号語を示す。2つのパリティチェック行は、互いに完全に独立して得られる。各情報シンボルは、各パリティ行に厳密に1度だけ現れる。すべてのパリティシンボルは、他の列（他のディスク）からの情報シンボルにのみ依存するものであって、互いに依存しあうものではない。したがって、1つの情報シンボルの更新は、2つのパリティシンボルのみを更新させる結果となる。

【0062】

上述のX符号は、真に対角線方向の計算を考慮して、素数nを使用している。しかし、nが素数でない場合にも、計算には別の線を使用できる。たとえば、n-1個のディスクのすべてを横断するような好適な所与のエンベロープも、X符号に従って使用することができる。これらの線はすべて、平行であることが好ましい。

【0063】

上述のように、X符号は列距離3を有し、2つの列消失または1つの列誤りを訂正することができるようにしている。消失とは、問題があつて、どの領域に問題があるのかがわかっている場合である。誤りは、問題の具体的な原因がわからない場合に生じる。復号化動作は、有限体演算を必要とせず、循環シフトおよび排他的ORのみを使用して、利用することができる。

【0064】

1つの消失の訂正は、どちらかのパリティ行を使用して傾き1または-1の対角線方向に沿って、その消失を簡単に復元することができる。

【0065】

サイズ $N \times n$ のアレイにおいて、2つの列が「消失」であるものと仮定する。この場合、これら2つの列の基本の未知のシンボルは、それらの列内にある情報シンボルである。各列が $(n-2)$ の情報シンボルを有するので、未知のシンボルの数は $2 \times (n-2)$ となる。同様に、残りのアレイは、これら $2 \times (n-2)$ の未知のシンボルを全て含む、 $2 \times n-2$ のパリティシンボルを含む。このため、消失の訂正は、 $2 \times (n-2)$ の線形方程式から $2 \times (n-2)$ の未知数を算出する問題となる。これらの線形方程式は、線形独立であるため、解くことが可能となる。

【0066】

さらに、この符号の同じ列内の2つの情報シンボルが、同じパリティシンボル内に現れることはない。したがって、各方程式は、2個以下の未知のシンボルを有することになり、未知のシンボルを1つしか含まない方程式もある。このことは、方程式を解く複雑さを大いに減じる。このシステムに従って使用されるシステムは、1つの知られている未知のシンボルを有する方程式から始める。これらの方程式を解くことは比較的簡単である。このプロセスは、すべての方程式が解けるまで、他の未知の解を求めて続けられる。

【0067】

消失列が i 番目および j 番目 ($0 \leq i < j \leq n-1$) の列であると仮定する。各対角線は、 $n-1$ 本の列のみを横断するので、1本の対角線が最終行において1本の列と交差する場合には、その列のどのシンボルも、この対角線内に含まれることはない。これにより、2つの消失列の1つのシンボルしか含まないパリティシンボルの位置が決定される。このシンボルは、この対角線に沿った簡単なチェックサムから復元することができる。

【0068】

まず、傾き1の対角線について考える。 i 番目の列の x 番目のシンボルが、1対角線内の唯一の未知のシンボルであると仮定する。この場合、この対角線は $(n-1)$ 番目の行において j 番目の列に当たり、 y 番目の列で第1のパリティ行に当たることになる。すなわち、3つの点 (x, i) 、 $(n-1, j)$ および $($

$n-2, y)$ は、傾き 1 の同じ対角線上にあり、したがって、以下の等式が成り立つ。

【0069】

【数3】

$$\begin{cases} (n-1)-x \equiv j-i \pmod{n} \\ (n-1)-x \equiv j-i \pmod{n} \end{cases}$$

$$(n-1)-(n-2) \equiv j-y \pmod{n}$$

式中、 $1 \leq j-i \leq n-1$ および $0 \leq j-1 \leq n-1$ であるため、 x および y の解は以下となる。

$$\begin{aligned} x &= \langle (n-1)-(j-i) \rangle_n = (n-1)-(j-i) \\ y &= \langle j-1 \rangle_n = j-1 \end{aligned}$$

【0070】

したがって、パリティシンボル $C_{n-2, j-1}$ は、 i 番目の列内のシンボル $C_{(n-1)-(j-1), i}$ の計算を可能にする。同様に、 j 番目の列内のシンボル $C_{(j-1)-1, j}$ は、パリティシンボル $C_{n-2, (j-1)-1}$ から直接解くことができる。

【0071】

傾き -1 の対角線方向と対称に、 i 番目の列内のシンボル $C_{(j-1)-1, i}$ はパリティシンボル $C_{n-1, (j+1)-n}$ から解くことができ、また、 j 番目の列内のシンボル $C_{(n-1)-(j-1), j}$ は、パリティシンボル $C_{n-1, i+1}$ から解くことができる。

【0072】

1つの情報シンボルと、傾き 1 および -1 の対角線とがそれぞれ厳密に 1 度だけ交差することに留意されたい。未知のシンボルが傾き 1 (または -1) の対角線に沿って解かれ、その後、その解かれたシンボルと交差する傾き -1 (または 1) の対角線に沿ってパリティシンボルが解かれると、他の列における別の未知のシンボルを解くことができる。この手順は、パリティシンボルが「消失」列となるか、または、解かれたシンボル自体がパリティシンボルとなるまで、反復的に用いることができる。これら同様の技術を使用して、1つ以上のどのような所

望の未知のシンボルも、復元することができる。

【0073】

好ましいシステムは、 $N = n$ または N として素数を使用する。図5および図6のようなシステム ($n = 4$; $k = 2$) を上述のように使用することもできる。

【0074】

分散型読出／書込

本システムは、分散型読出および書込システムを使用することによって、新しい種類の動作を行なうことができる。

【0075】

情報の冗長記憶により、システムは、システムの帯域幅を最大にするように、 n 個のノードすべてから読出しを行なうことができる。この場合、システムは、ノードの生の情報部分 1502 からのみ読出しを行なう。

【0076】

代替的に、ノードのうち k 個のみを読出すこともできるが、これら k は、それらのパリティ部分 1504 とともに読出される。従来の「訂正」とは異なり、このシステムは、利用可能なクラスタのうちどれを使用するかを、システムが見たネットワークの状態に基づいて選択する。異なる部分を、たとえば奇偶符号等の異なる符号のために使用することができる。

【0077】

分散型の書込は、情報が変化するたびに、影響を受けるすべてのノードに書込みを行なう。しかし、更新はできる限り小規模とされる。MDS 符号は冗長性を保証し、更新を最適に最小かつ効率的にする。平均単位パリティ更新数は、符号内で単一の情報ビットが変化した場合に影響を受ける、パリティビットの平均数を表わす。このパラメータは、記憶アプリケーションにアレイ符号が使用されているときに特に重要となる。X 符号は、各単一の情報ビットの変化が2つのパリティビットの更新しか必要としないという意味で、最適である。

【0078】

X 符号の別の重要な特徴は、独立なパリティビットが形成されることに起因する。これまで使用されてきた符号の多くは、3 という符号距離を形成するために

従属パリティ列に依存している。これらのパリティの計算は、パリティが互いに依存しあっているために、非常に複雑であり得る。このことはしばしば、その符号の平均単位パリティ更新数がアレイの列の数に応じて線形に増加する、という状況を生み出す。

【0079】

米国特許番号第5, 579, 475号に記載されたEVENODD（奇偶）符号または他の同様のシステムは、独立パリティ列を使用して、情報の更新をより効率的なものにしている。

【0080】

検出

分散型データ記憶システムは、複数ノードにわたってサーバ機能を広げる。これは、本システムに従えば、多数ノードの各々の上で稼動する、特別な通信層を使用して行なわれる。この通信層は、アプリケーションに対して透過的である。特別な分散型読出システムおよび分散型書込システムは、本システムの頑健な動作もまた維持する。

【0081】

好ましいシステムの通信アーキテクチャを図7に示す。実際の通信およびネットワークインターフェイスは、要素700として示されている。通信は、イーサネット、ミリネット、ATMサーバネット（ATM Servernet）、または従来の他のどのような通信方式をも含む、従来の方法で行なうことができる。これらの従来のネットワークインターフェイスは、冗長通信層によって制御される。

【0082】

通信は、ネット監視（「NETM」）プロトコルシステム702によって監視される。NETMは、各ノードにおいてチャネル状態およびチャネル状態の履歴を判断する接続プロトコルを維持する。より特定的には、NETMは、NETMがその上で稼動しているローカルノードから各リモートノードへの、そのローカルノードとリモートノードとを結ぶ各接続経路を介した接続のすべてを監視する。NETMは接続チャートを保持するが、このチャートは、ローカルノードから各リモートノードへの可能なすべての接続のステータスを常時知らせる表示を含

む。

【0083】

実際の通信は、高信頼ユーザデータプロトコル（「RUDP」）によって制御される。RUDPは、ローカルノード（「ノードA」）が他のノード（「ノードB」）への通信を求めるリクエストに基づいて動作する。RUDPは、NETMから、ノードAからノードBへの正常に動作している通信経路に関する接続情報を入手する。RUDPは、NETMによって集められた情報を使用して通信経路を選択し、その情報を束ねられたインターフェイスを使用して送出する。RUDPはまた、よく知られているプロトコルシステムを使用して情報を適宜パッケージして、順序正しく確認された配信を行なう。

【0084】

NETMシステムは、システムの各ノード上で稼動して、システムに関する情報を見つける。NETMは、自身がその上で稼動しているノードを、ローカルノードとして見ている。NETMはシステムクルーを使用して、そのローカルノードとシステム内の他のすべてのノードとの間の接続の状態を判断する。

【0085】

すべてのノード上で同じプロトコルが稼動しているので、各ノード上の各NETMプロセスは、AからBへのどの接続状態についても、同じ条件を判断するであろう。NETMはまた、履歴チェック機構を使用して、すべてのノードがある時間にわたってチャンネル状態の同じ履歴を見るようにしている。

【0086】

好ましいシステムクルーは、ノードAからシステム内の他の各ノードへと、そのノードへの可能な各経路上を送られるメッセージから得られる。これらのメッセージは「ハートビート」と呼ばれる。NETMは、ローカルノード（「ノードA」）から各リモートノード（「ノードB」）へと、各経路上でメッセージを送る。各接続は、 i = ローカルインターフェイス、 j = リモートノード、および k = リモートインターフェイスを含む、 $C_{i,j,k}$ 「タプル」と呼ばれる情報の3つの項目によって特徴付けられる。このタプルは、明確な経路を規定する。

【0087】

NETMは、ハートビートを使用して、各経路 $C_{i,j,k}$ 上にAとBとの間に運用中の通信リンクが存在するかどうかを判断する。NETMプロトコルはノードB上でも稼働しているので、そのリモートNETMも、各経路 $C_{i,j,k}$ を介したノードBからAへの接続状態について、同じ判断を行なうであろう。

【0088】

たとえばバッファオーバーフロー等のような障害は、一方方向のみでチャンネルが失われる原因となり得る。接続プロトコルはトークンパッシングシステムを使用して、チャンネルの履歴を対称にする。

【0089】

履歴の検出は、接続の動作可能性に関するヒントの集積に基づく。ハートビートは好ましいヒントであり、ここにさらに詳細に説明する。別のヒントの例としては、通信ハードウェア、たとえばミリネットカードからの、障害表示がある。経路X上の通信を制御しているミリネットカードが、動作不能であるとの表示をすると、プロトコルはその経路が動作不能であると評価することができる。

【0090】

ヒントの集積は、Xを介したAからBへの通信経路の状態を評価する変数の状態を設定するのに使用される。この変数は、稼働中（アップ）に対してU、ダウンに対してDの値を有する。

【0091】

この動作を、図8の概略的なフローチャートに示す。図8の実施例は、信頼性の低いメッセージで形成されたハートビートメッセージを使用している。高信頼メッセージ通信システムでは、送信側のノードがメッセージの受信確認を受取ることが要求される。送信側ノードは、何らかのメッセージ受信確認を受取るまで、そのメッセージを送り続けるのである。これに対し、図8のシステムは、信頼性の低いメッセージ通信を採用している。すなわち、メッセージが単に送信されるのみであり、受信の確認が受取られることはない。

【0092】

メッセージ800は、低信頼パッケージメッセージとしてノードBに送られる。ハートビートは、好ましくは10ミリセカンド（ms）毎に送られる。システ

ムは、ステップ802においてネットワークヒントを待ってチェックして、ネットワークリンクの状態および履歴を評価する。ハートビートは、あるノードから他ノードに送信されるどのようなメッセージであってもよい。

【0093】

各ノード上では同じプロトコルが稼動しているので、各ノードは、10ms毎に他の各ノードからハートビートを受信すべきであることをわかっている。各NETMにおいて、そのNETMが他のノードからハートビートを受信するたびにリセットされるタイマが稼動している。他のノードからハートビートを受信する前にタイマが時間切れとなった場合、その接続に何か問題がある、と判定することができる。

【0094】

各サイドは、ある時間にわたって同じ履歴を見ていることを確認しようとする。これは、ポイント・ツー・ポイントプロトコルを構成しているノードの対の間で高信頼トークンを渡すことによって行なわれる。各トークンは、ノードがイベントを見たことを示す。そのトークンが他のノードによって受取られる場合、そのノードもまた、対応するイベントを見てトークンを送ったことを意味する。各サイドは、イベントを認めるとトークンを渡す。これにより、両サイドにおける履歴を同じものに維持することができる。

【0095】

各サイドが渡すことのできるトークンの数は有限である。これは、他のノードによってイベントが確認される前に起こり得るイベントの数を制限する効果がある。たとえば、1サイド当り2つのトークンが当初あったとすると、ノードは渡せるトークンを2つしか持っていないことになる。チャネルの状態に変化が認められるたびにトークンが1つ渡される。他サイドからトークンが届かなければ、2つの変化を認めた後には、そのノードの所有するトークンはなくなる。このことは、各ノードが、他ノードよりも2つのイベントまたはアクションのみしか先行する（または遅れる）ことができないことを意味する。このトークンの手渡しは、2つのノード間の自由度を制限する。すなわち、1つのノードがチャネルがダウン状態であるという報告を受けてから、他ノード側がその情報を得るのを待

つ間、それら2つのノードが離れて動作することのできる度合いを制限する。

【0096】

別の見方をすれば、それらトークンは、1つのノードが他ノードが同様にアクションを取ったことを知る前に行なうことのできる、遷移の最大数を設定する。

【0097】

NETMシステムの好ましい実施例を、図9の接続プロトコルステートマシンにおいて示す。図10Aおよび図10Bは、その実施例のフローチャートである。ステップ1000は、1つのノードから他の知られているすべてのノードへの物理的に可能なすべてのチャンネルの各々について、ローカルインターフェイスID、リモートマシンID、およびリモートインターフェイスIDを含む $C_{i,j,k}$ の3タプルを形成する最初のステップを含む。プロセス $ConnP(C_{i,j,k})$ は、すべての $C_{i,j,k}$ の3タプルに対して実行されて、これらチャンネルの各々について接続状態が判断される。これにより、各 C_i チャンネルについてアップ/ダウン(1または0)のステータスを示すブール値を記憶する、 $Connected(C_{i,j,k})$ と呼ばれるデータ構造が形成される。

【0098】

ステップ1002は、 $ConnP(C_{i,j,k})$ イベントがあったかどうかを判断する。もしイベントがなかった場合、実行すべきことは何もないので、プロセスは初めに戻る。

【0099】

ステップ1002においてイベントが検出されれば、フローはステップ1004に進み、そこで、そのイベントがシステムアップイベントであるかどうか判断される。そうであれば、結果として「1」が戻され、そうでなければ、結果として「0」が戻される。

【0100】

図10Bのリンクステータスのフローチャートは、他のエンドポイントシステムの動作の証拠として「トークン」のカウントを使用する。

【0101】

ステップ1010で、プロセスは、最初の値 $n \geq 2$ に設定されているトークン

カウント（「t」）で開始する。システムは、ステップ1012で、アップ（「1」）に最初に設定されている状態で開始する。ステップ1014は、タイム・インイベントがあったかどうかを検出する。タイム・インイベントは、たとえば、ノードBからハートビートを受信することによって引起される。状態はこの時点で既にアップなので、タイム・インイベントが検出されても状態はアップのままであって、さらなるアクションが取られることはない。ステップ1014においてタイム・インイベントがなかった場合には、ステップ1016において、たとえばタイマが時間切れとなる前に予測されるハートビートが受取られなかったために生じる、タイム・アウトイベントの判断が行なわれ、ステップ1018においては、トークンが受信されたかどうか（「トークン到着イベント」）が判断される。それらのイベントが起こっていなかった場合には、制御は再びステップ1012に渡され、それらのイベントが起きないかどうかの監視が続けられる。システムはこの時点では必ずトークンを持っているため、別のトークンを調べる必要はない。

【0102】

ステップ1016においてタイム・アウトイベントがあるということは、経路Xを介してノードBからハートビートが受信されなかったことを意味し、したがって、経路Xを介したノードBへの通信に問題がある可能性を意味する。そこで、制御がステップ1020に移る。ステップ1020においては、タイム・アウトイベントを示すトークンがノードBに送られて、所定の時間内にハートビートが届かなかったことが報告される。トークンが送出されたので、トークンのカウントもまた1020においてデクリメントされる。この後、ステップ1022においてConnPの状態がDに変更される。

【0103】

ステップ1018においてトークン到着イベントがあると、1024においてトークンの受信ステップが行なわれ、トークンのカウントがインクリメントされる。現時点におけるトークンカウントが1026における最大トークン値nよりも小さい場合には、トークンカウントは1028においてインクリメントされる。トークンが失われることがあるので、他方端における遷移は、トークンパッシ

ング方式によって許容される許容可能な自由度の範囲内であり、受取られたトークンによって、両サイドが同期状態に戻される。

【0104】

トークンカウントがN未満でない場合には、トークンカウントはその最大値である。したがって、システムは遷移をする必要がある。これは、1030においてトークンを送出することによって行なわれ、その後、システムがダウンする。これは、 $ConnP \rightarrow 0$ または、1022においてDとして示される。これにより、ダウンルーチン処理動作が開始される。

【0105】

ダウンルーチン処理動作は、アップルーチン処理動作と同様である。タイム・アウトイベントが1030で検出される。これは、システムが既にダウンしているので何ら影響を及ぼすことはない。タイム・インイベントが1032で検出される。このタイム・インイベントは、もしその遷移を示すよう送ることのできるトークンがあれば、システムをアップ状態に戻すことができる。このルーチンは、ステップ1040においてトークンを調べる。利用可能なトークンがなければ、遷移は起こらず、フローは1022に戻る。手渡せるトークンがあれば、フローは1042に進み、トークンカウントがデクリメントされる。 $ConnP$ の変数は、そのアップ状態に戻り、トークン処理ルーチンが始まる。

【0106】

経路Xを介したノードAからノードBへの各システムは、NETMプロトコルによってこのように特徴付けられる。

【0107】

アプリケーションは、RUDPの上で実行される。たとえば、プロセスIDを有するアプリケーションは、最初に、システムに対して自身の身元を明かす。たとえば、このアプリケーションは、プロセス6として自身を識別しかつプロセス4に対して送信の希望を示す、メッセージを送信することができる。このアイデンティフィケーションは、上述の Ci, j, k タプルを使用する。NETMは、このオペレーションのための通信経路を決定する。

【0108】

実際の通信は、一旦決定されると、いわゆるスライディングウィンドウプロトコルを使用して行なわれる。スライディングウィンドウはよく知られており、たとえば、米国特許番号第5,307,351号に記載されている。スライディングウィンドウは、データパケットを適切にパッケージングすることによって、高信頼メッセージ通信方式を監督する。スライディングウィンドウは、実質的に、シーケンス番号および確認を管理する。データは高信頼パケットとして送出され、2つ以上のウィンドウが送出される前に受信側が受信を確認することが要求される。受信が正しく確認されると、情報のウィンドウは次の、まだ確認されていない情報のパケットに「スライド」する。

【0109】

RUDPは、スライディングウィンドウモジュールを使用して実際の通信を行なう。RUDPはまた、有効な情報経路を提供するよう、NETMをコールする。ノード間に使用できる経路が2つ以上ある場合には、RUDPはそれら使用可能な経路を巡回する。

【0110】

RUDPはまた、NETMによって提供される情報を使用してシステムを再設定することによって、論理ネットワークとしても機能する。

【0111】

基本的なRUDPのフローチャートを図11に示す。オペレーションは、ステップ1100において受信イベントを判断することから開始する。ステップ1100において受信イベントが受信されなければ、ステップ1102において、送信イベントがあったかどうか判断される。もし送信イベントがなければ、LNETは何もすることがなく、フローは戻って、イベントのチェックを続ける。

【0112】

ステップ1100において受信イベントが検出されれば、フローはステップ1110に進む。ステップ1110は、データが何らかの $C_{i,j,k}$ タブルを示しているかどうかを判断する。もしそうでなければ、ステップ1112において誤りが判断される。

【0113】

正しいデータが得られれば、そのデータはステップ1114において受信されて、ステップ1116においてシステムに戻る。

【0114】

送信イベントは、送出されるデータおよびイベントを受信するリモートマシンを示す、 $C_{i,j,k}$ の引数を要求する。このためには、リモートマシンに対して動作の引数の1つとして示されるアップチャネル $C_{i,j,k}$ が存在するかどうかを、ステップ1120で判断せねばならない。もし存在しなければ、ステップ1122において接続損失誤りが宣言されるが、より一般的なケースとして、少なくとも1つのアップチャネルが存在すれば、そのアドレスは $C_{i,j,k}$ タプルの引数を使用している。その後、プロセスは1130に戻る。

【0115】

プロセス1120は、NETMを使用して、ローカルマシンからリモートマシンへの既存の経路を調べる。したがって、NETMは、LNETがデータ構造を使用する間も、データ構造を維持する。

【0116】

情報サーバ

ここに記載するシステムは、情報サーバ、すなわち、リクエストに応じてユーザに情報を提供するサーバ、に特別に応用される。情報サーバは、インターネット（ウェブ）サーバ、ビデオサーバ、または、情報を提供する他のどのような種類のデバイスであってもよい。

【0117】

システムは、どのノードも、他のいずれかのノードまたはノードの組合せからそこに記憶されているどのような情報もリクエストすることができる、という意味で、サーバとして使用される。たとえば、25個の異なるノードから情報を要求するリクエストを行なうこともできる。このシステムは、最も近くにある25個のノード、または、最も使用頻度の低い25個のノードを選択することができる。これにより、システムは、過負荷のノードを、それらがあたかも障害を起こしているかのように無視することができる。

【0118】

システムがビデオサーバとして使用される場合、送出されるべきビデオの記憶場所としては、システム上のあらゆる場所が可能である。本方式に従えば、ビデオは分散情報として、ネットワークの種々のノード間で記憶されており、特定のネットワークが故障した場合にもビデオ情報が取出せるようにしている。

【0119】

サーバシステムは、ビデオを記憶しているノードから、そのビデオが提供されるようリクエストする。このシステムの特別な技術は、特定数の故障によってシステム全体の動作が決して中断されることのないようにする。たとえば、2つのノードが故障しても、記憶情報を得ることはできる。というのも、情報はネットワーク内の他のローケーションに冗長に記憶されているためである。

【0120】

別の応用例としてウェブサーバがある。ウェブサーバは、TCP/IPプロトコルおよびパケット通信を使用して、インターネットの情報を取得する。やはり、この情報も、分散型サーバ内のどこに記憶されていてもよい。通信または記憶を問わず、どのような2つの障害も、情報の取得を妨げることはない。

【0121】

このシステムは、他にも、拡張および修復に応用される。どのノードもどの時点においても取除くことができ、システムの残りの部分はそれでも中断なく動作し続ける。そのノードは、空白のノードに置換されてもよく、その場合、ネットワークは、それが空白の列であると認める列に、冗長データを利用して情報を書き込み始める。

【0122】

以上に、いくつかの実施例のみを詳細に開示したが、当業者には、開示された実施例の範囲内に他の実施例が存在すること、また、開示した実施例から、本発明を実施するための技術が他にも想定され得ること、が理解されるであろう。

【図面の簡単な説明】

【図1】 最も簡単なネットワークの例を示す、基本的なブロック図である。

【図2】 より多くのスイッチおよびより多くの計算ノードを有する、より

複雑な例を示す。

【図3】 より信頼性の高いネットワーキングの例を示す。

【図4】 フォールトトレラントシステムを示す。

【図5】 このシステムを使用してビデオを記憶する方法の例を示す。

【図6】 システムがどのようにリンク障害を許容することができるかを示す。

【図7】 信頼性の高い通信のための、ソフトウェアアーキテクチャのブロック図である。

【図8】 ネットワーク監視プロセスのための基本的なソフトウェアフローチャートを示す。

【図9】 ネットワーク監視プロセスのための接続プロトコルステートマシンを示す。

【図10A】 接続のためのデータ構造のフォーメーションを示す。

【図10B】 リンクステータスオペレーションのフローチャートを示す。

【図11】 RUDPプロセスのフローチャートを示す。

【図12】 計算ノードおよびスイッチング要素の可能な配列を示す。

【図13】 計算ノードおよびスイッチング要素のより高度な配列を示す。

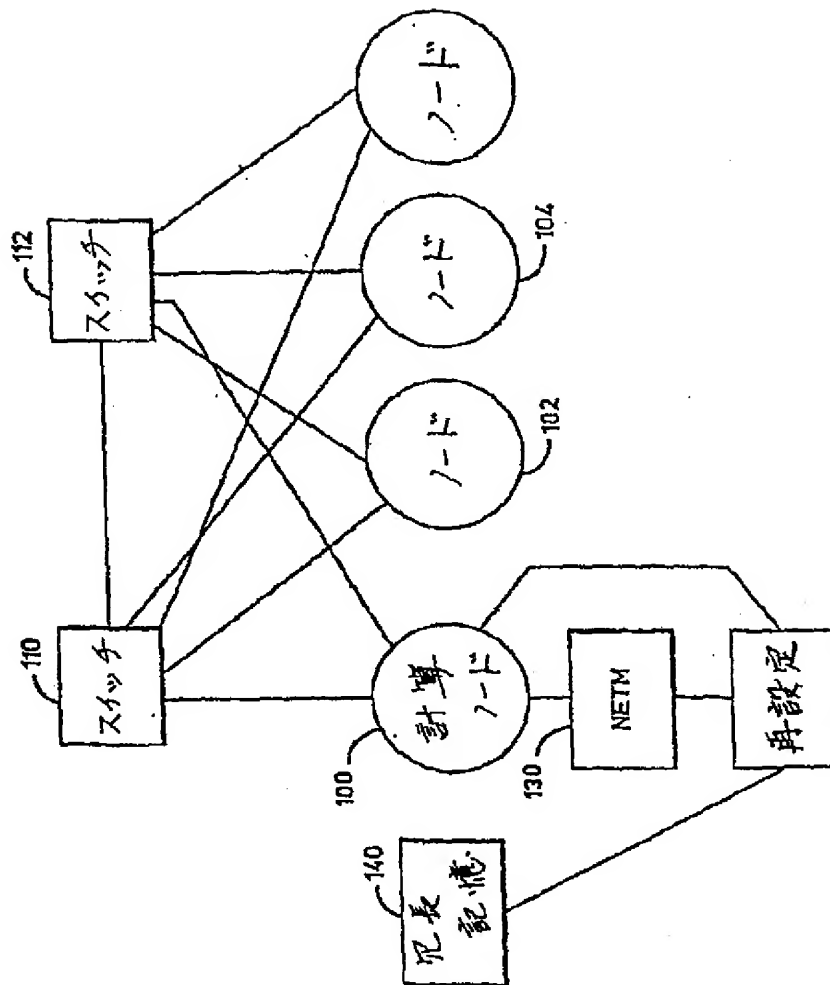
【図14】 X符号における、5×5のアレイ符号のための、パリティ行の計算を示す。

【図15】 X符号システムの基本的なレイアウトを示す。

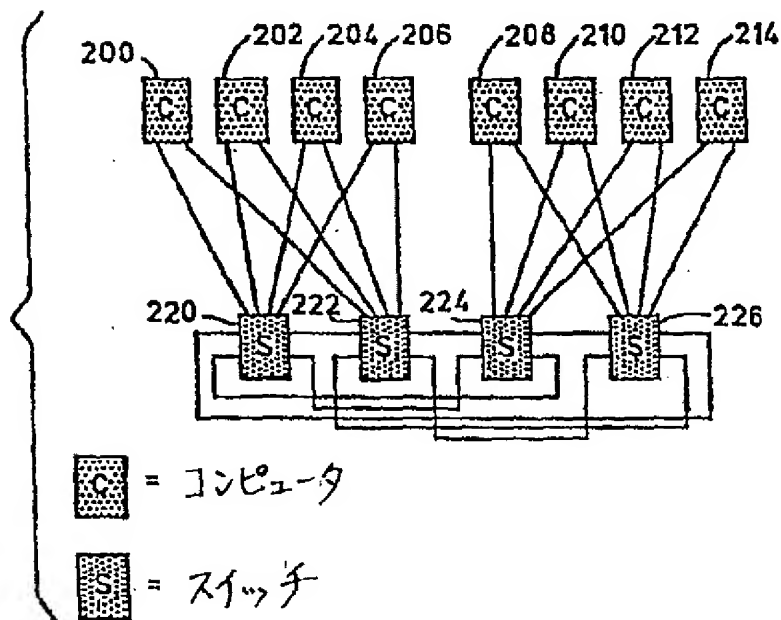
【符号の説明】

100, 102, 104, 106 計算ノード、110, 112 スイッチ、
130 ネットワーク監視 (NETM)、140 冗長記憶。

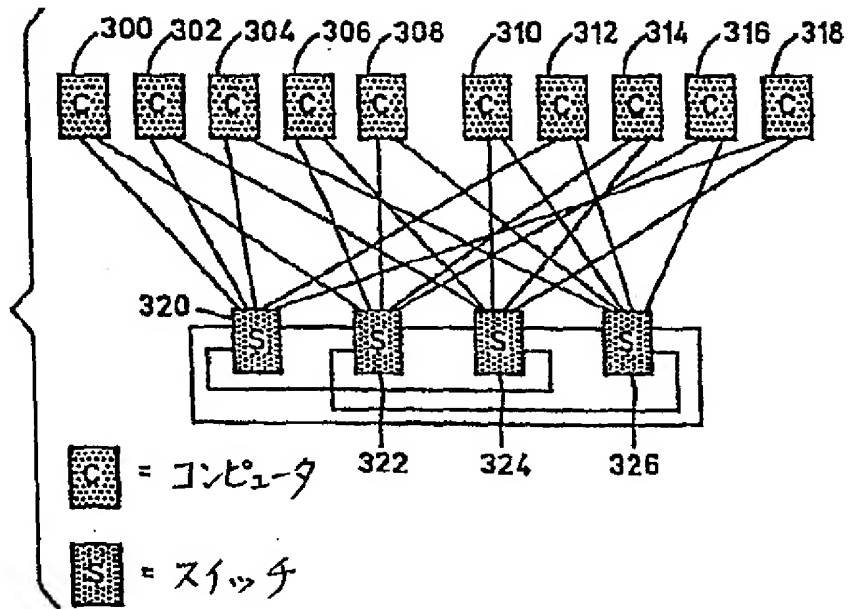
【図1】



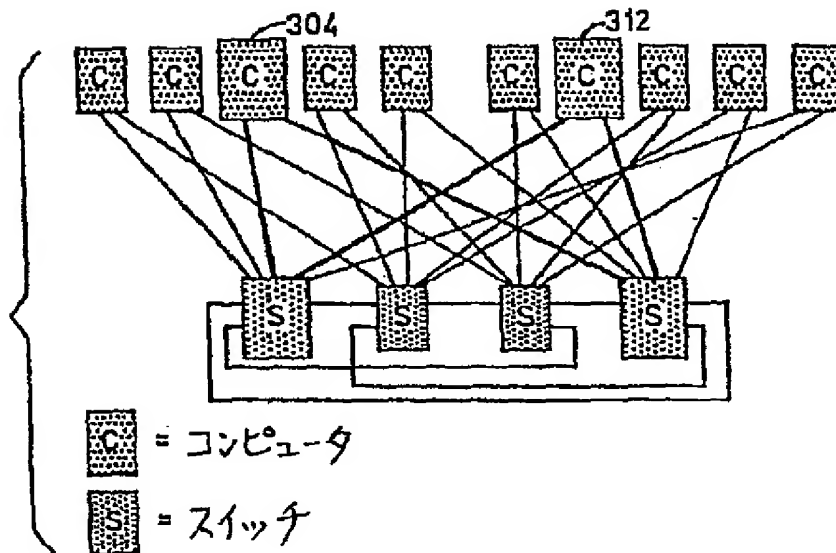
【図2】



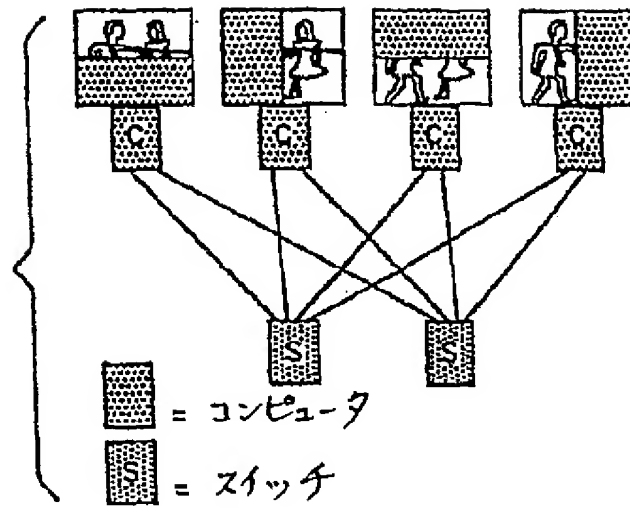
【図3】



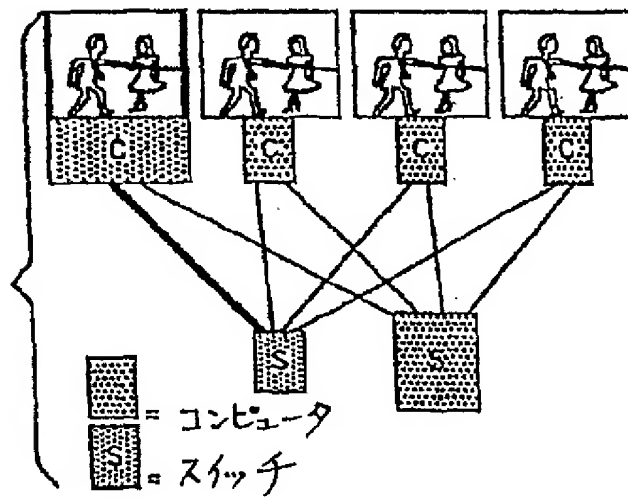
【図4】



【図5】



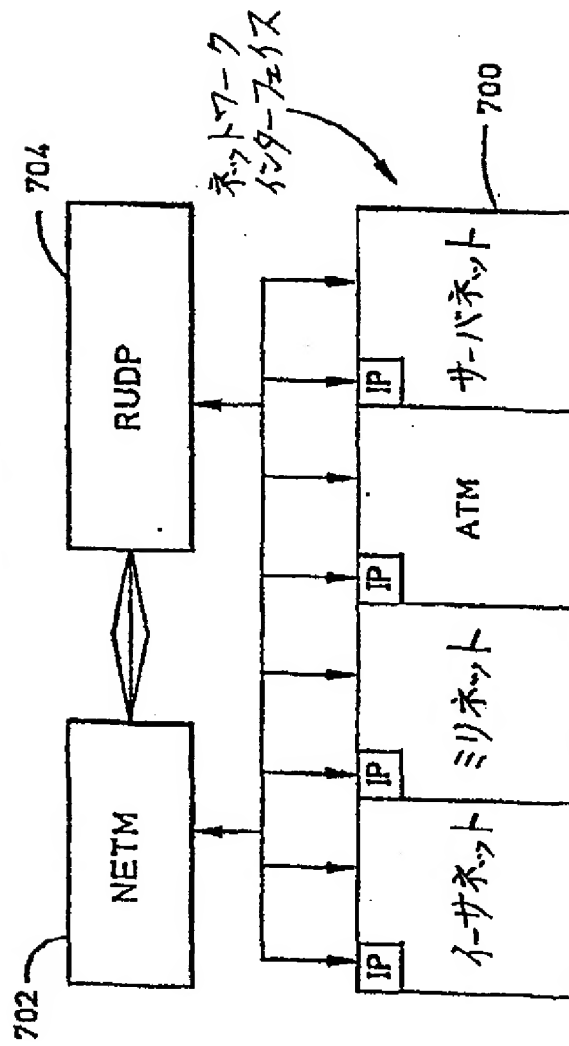
【図6】



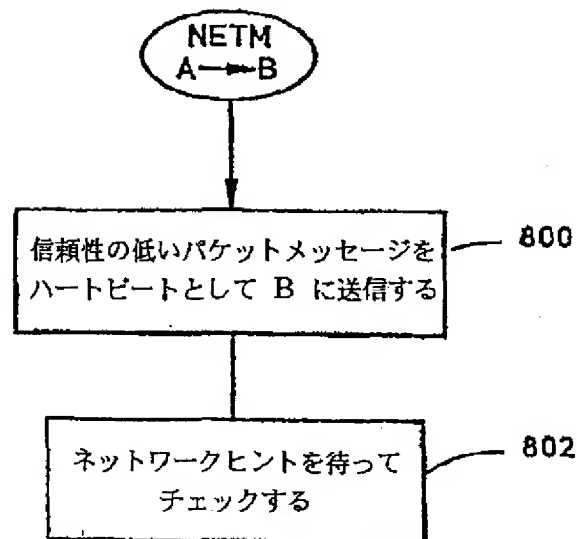
【図7】

高信頼通信のためのソフトウェアアーキテクチャ

NETM(ネットモニタ): チャネル状態の適時検出および、各エンドポイントにおけるチャネル状態が同一履歴となることを保証するために、接続プロトコルを使用する。
 RU DP(高信頼ユーザデータプロトコル): 順序正しく確認された配信および束ねられたインターフェイスを併せて提供し、論理ネットワークとする。

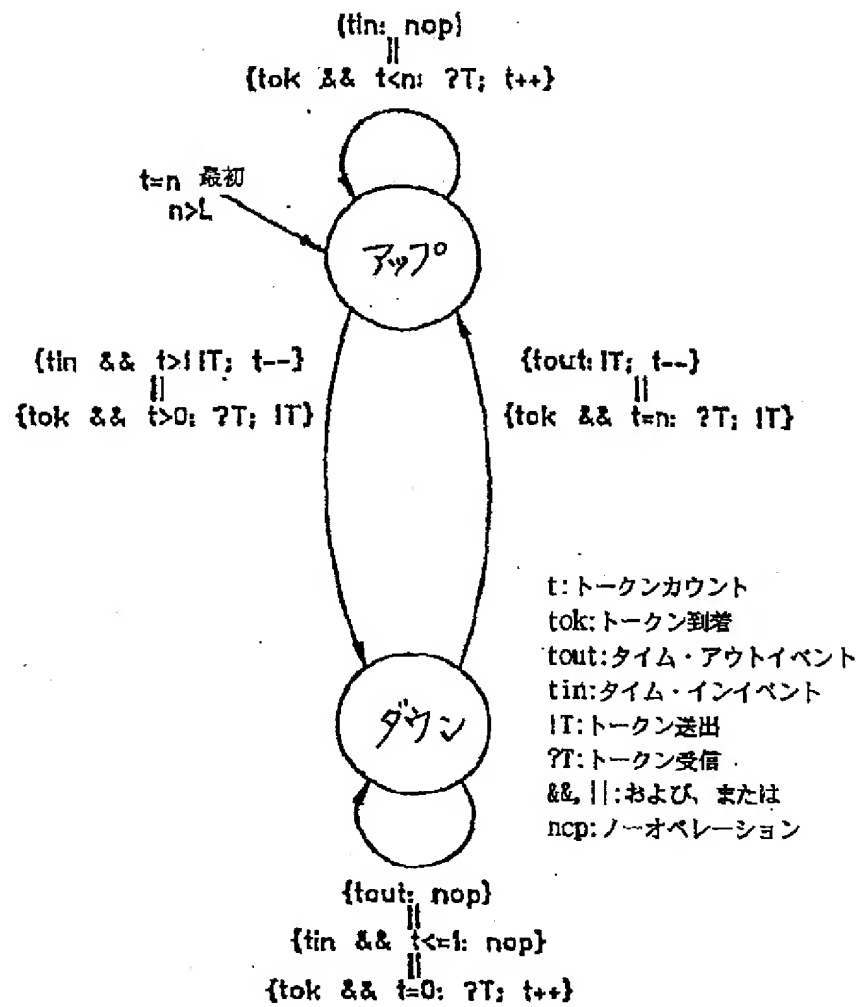


【図8】

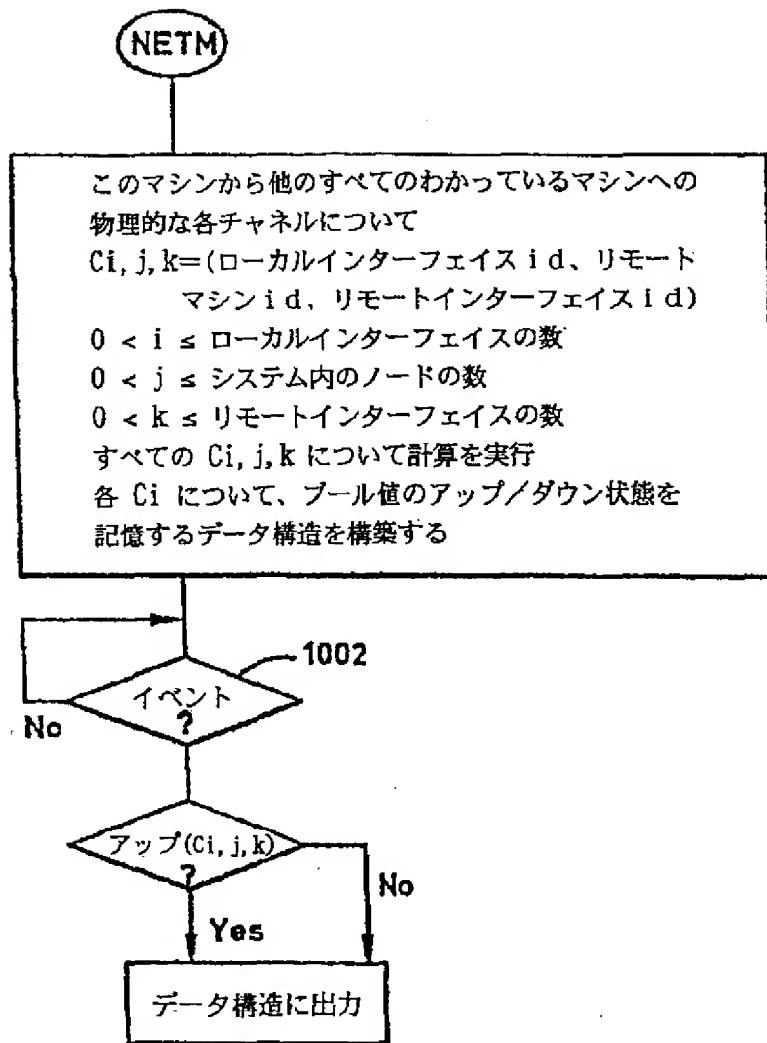


【図9】

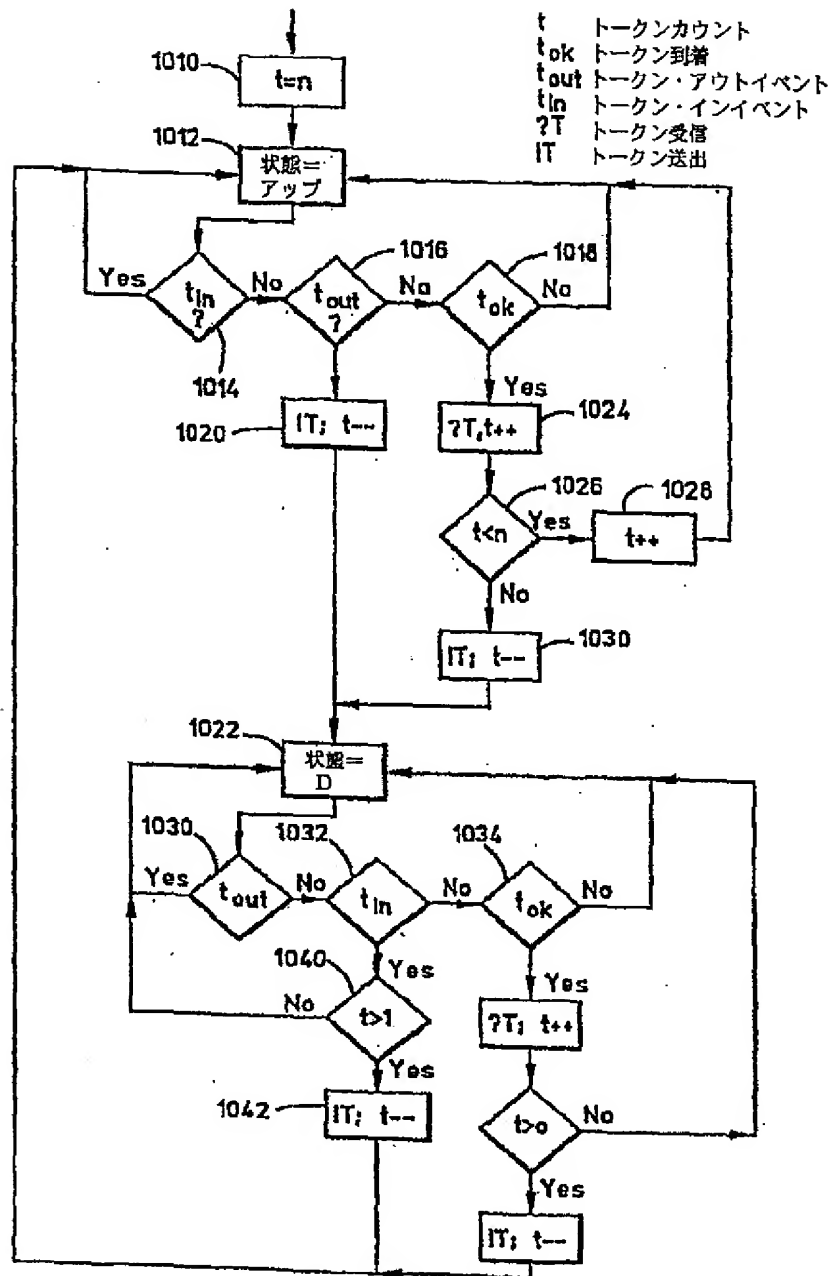
接続プロトコルステートマシン



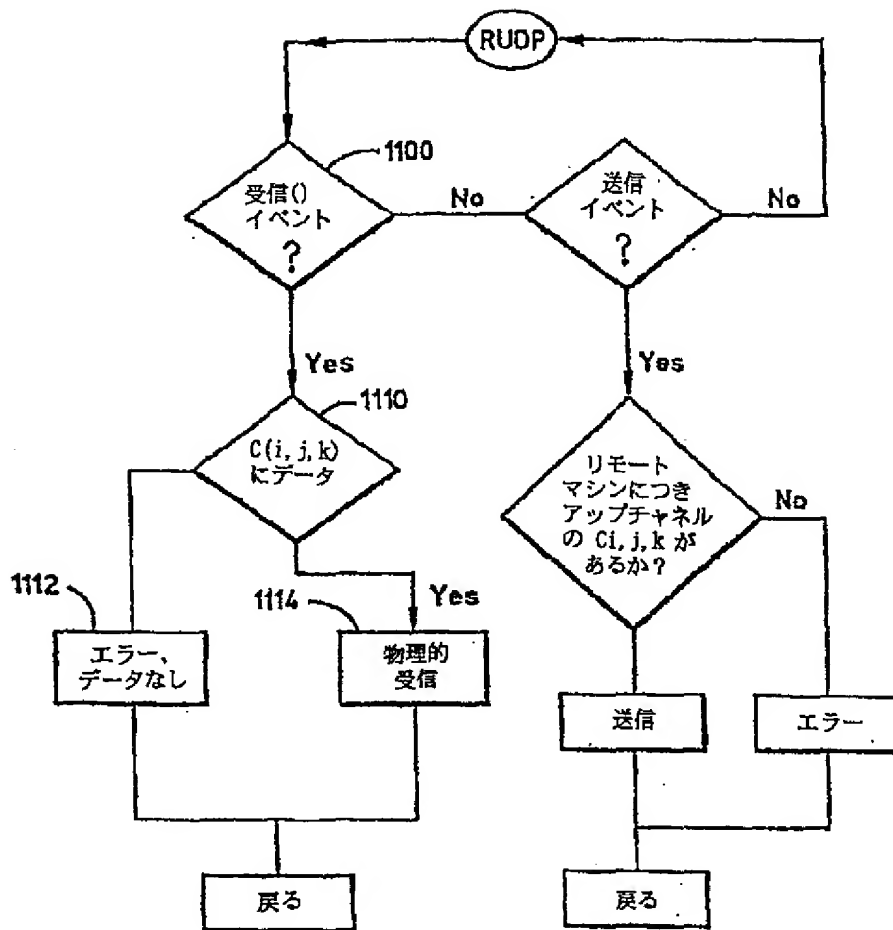
【図10A】



【図10B】



【図11】



【図12】

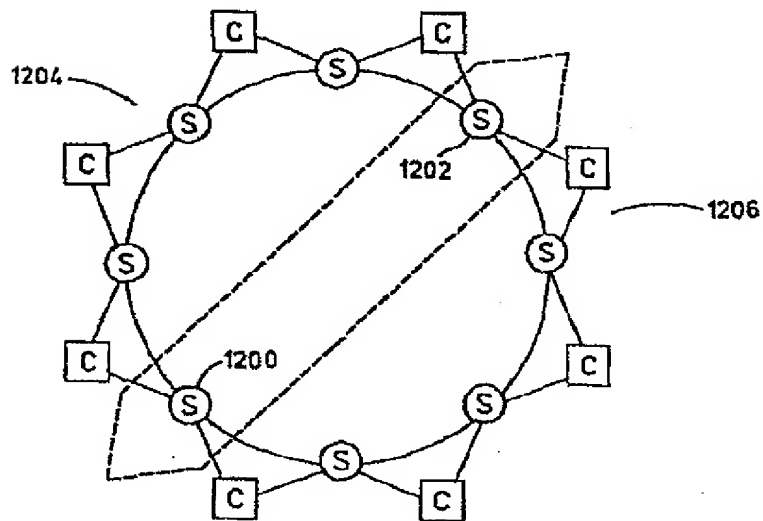
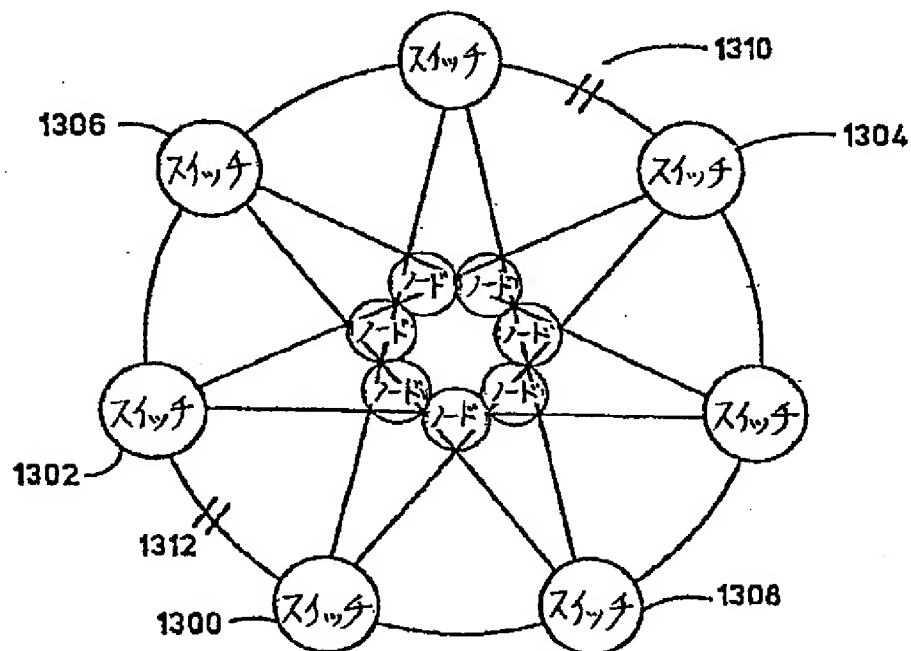


FIG. 12

【図13】

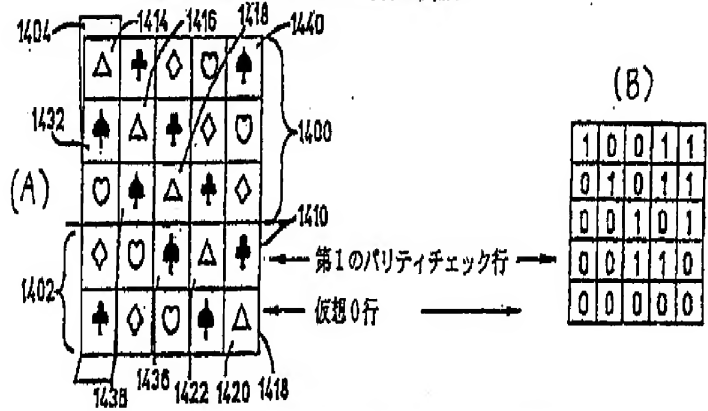
| |
|------------------------------|
| 一般的な問題 |
| 設定：スイッチおよびノードのネットワーク。 |
| 目標：ノード間の通信。 |
| 障害：スイッチ、ノードまたはリンクの故障。 |
| 具体的な問題 |
| 設定：スイッチはパケットを前方に送り、ノードは送らない。 |
| 目標：分離されるノードの数を一定に保つ。 |
| 障害：スイッチの故障。 |



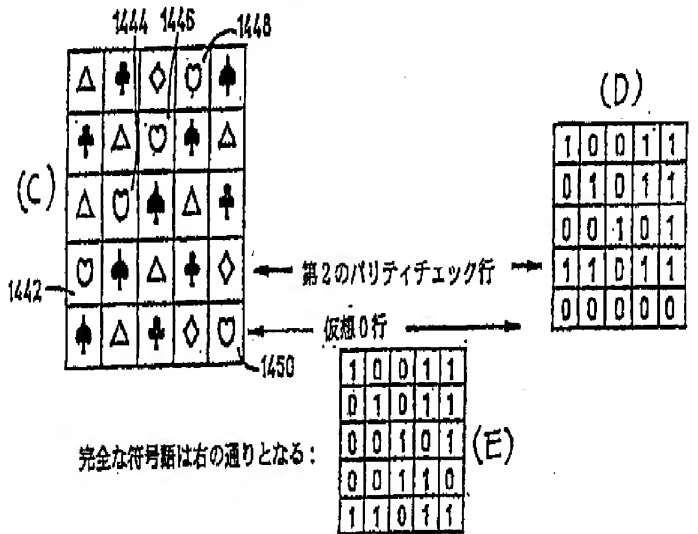
例1: 5×5のアレイ符号

第1のバリティ行は傾き1の対角線に沿って計算される。

最終行は、仮想ゼロ行である。以下を参照:

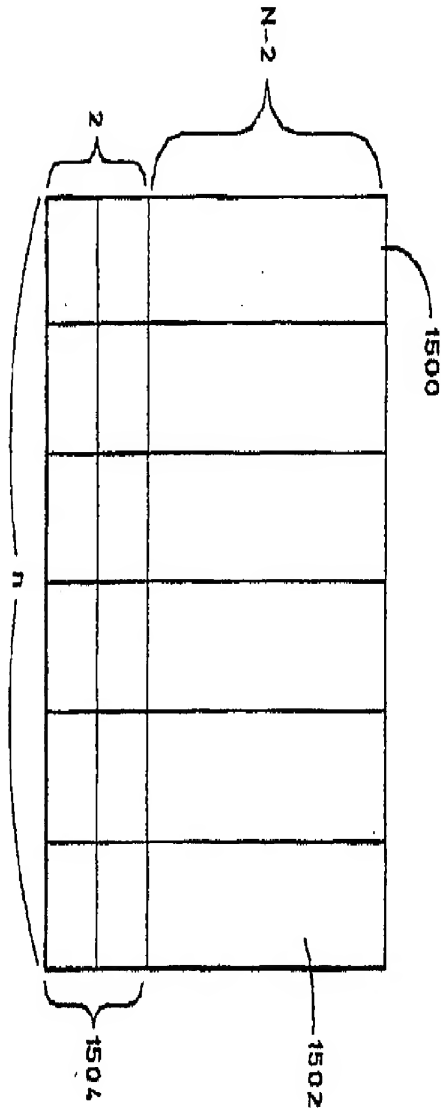


第2のバリティ行は傾き-1の対角線に沿って計算される。以下を参照:



完全な符号語は右の通りとなる:

FIG. 15



【手続補正書】特許協力条約第34条補正の翻訳文提出書

【提出日】平成12年3月31日(2000.3.31)

【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】特許請求の範囲

【補正方法】変更

【補正内容】

【特許請求の範囲】

【請求項1】 冗長分散型ネットワークシステムであって、

複数のシステムノード(100, 102, 104, 106)を含み、前記システムノードの各々は少なくとも2つの通信装置および記憶装置を有し、前記記憶装置はいずれか1つのノードにおける記憶情報が他のいずれかのノードにおけるデータと組合されたときに記憶されたデータの再構築を可能にする、ネットワークのための情報の冗長記憶を有し、さらに

前記システムノードの前記通信装置に、いずれか1つのシステムノード内の前記通信装置の各々が複数のスイッチング装置の異なる1つに接続されて、前記システムノードの各々が少なくとも2つの異なる経路のうちの1つを介して互いに通信することができ、したがって冗長通信を提供するように接続される、複数のスイッチング装置(110, 112)と、

ネットワークシステムの動作を妨害し得るシステム機能状態を検出する、各ノード(100, 102, 104, 106)における検出ルーチンと、

必要に応じてネットワークの冗長性を利用してシステム機能状態を補償するよう、前記通信冗長性または前記記憶冗長性のうち少なくとも1つを使用してネットワークを再設定する、論理ネットワークプロセス(140)とを含む、システム。

【請求項2】 前記検出ルーチンは、前記システムノードの少なくとも複数上の機能状態を検出し、前記システムノードの各々上では各前記検出ルーチンによって同じプロトコルが実行される、請求項1に記載のシステム。

【請求項3】 前記システム機能状態は、ネットワーク通信内の障害、メモ

り記憶内の障害、または、望ましくない結果をもたらす他の種類の障害を含む、請求項1に記載のシステム。

【請求項4】 前記システム機能状態は、ネットワーク通信における障害を含み、前記論理ネットワークプロセスは、接続が異なる接続に変更されるよう指令する、請求項3に記載のシステム。

【請求項5】 前記システム機能状態がメモリ記憶における障害を含むときには、前記論理ネットワークプロセスは、所望の情報が前記冗長データ記憶から得られるよう指令する、請求項4に記載のシステム。

【請求項6】 前記検出ルーチン(702)は各ノード(100, 102, 104, 106)で動作して、ネットワーク内の他のノードに対する接続状態を見る、請求項1に記載のシステム。

【請求項7】 前記検出ルーチン(702)は、ネットワークの状態に関するヒントを使用して前記接続状態を判断するよう動作する、請求項6に記載のシステム。

【請求項8】 前記ヒントは、所定の間隔で各前記ノードによって生成されるハートビート信号を含み、前記検出ルーチンは、前記ハートビート信号を受信して前記ヒントの1つとしての前記ハートビート信号の存在または不在を検出するよう動作する、請求項7に記載のシステム。

【請求項9】 トークンパッシングシステムをさらに含み、各ノード(100, 102, 104, 106)は監視されるチャネルを介して監視されるノードにおけるイベントを判断し、前記イベントを示すよう前記監視されるチャネルを介して前記監視されるノードにトークンを渡し、前記監視されるノードは、前記イベントに基づいた動作を示すよう前記トークンを戻し、各ノードは、該ノードとは異なるノード上で対応するイベントが起きずとも該ノードにおいて起こすことのできるイベントの数を制限するよう、所定数のトークンのみを有する、請求項7に記載のシステム。

【請求項10】 各前記ノードがネットワークの同じ履歴を見ることを保証する手段をさらに含む、請求項6に記載のシステム。

【請求項11】 前記接続は、計算ノードのグループが分離されることのない

いように行なわれる、請求項1に記載のシステム。

【請求項12】 前記スイッチ (110, 112) は、可能な限り最も非ローカルな方法で前記ノードを接続する、請求項1に記載のシステム。

【請求項13】 前記スイッチは、互いに最も遠い2つのノード間を接続する、請求項12に記載のシステム。

【請求項14】 前記接続は、どの2つのノードの故障によってもノードの1グループをノードの他のグループとの通信から分離することのないように行なわれる、請求項13に記載のシステム。

【請求項15】 各ノード (100, 102, 104, 106) は、少なくとも2つの経路によって各他ノードと接続され、

ネットワークモニタ (702) をさらに含み、該ネットワークモニタは、各前記ノードにおいて稼動して、前記ネットワークモニタがその上で稼動するローカルノードから各リモートノードへの各接続経路上のすべての接続を監視する、請求項1に記載のシステム。

【請求項16】 高信頼ユーザデータプロトコルをさらに含み、該プロトコルは、前記ローカルノード上で稼動し、該ローカルノードから他のノードへの通信リクエストを受信し、かつ、前記ネットワーク監視プロセスから経路を判断する、請求項15に記載のシステム。

【請求項17】 ノード間の物理的な接続を、異なるノード接続に変更することを可能にする、論理ネットワーク相互接続を使用した前記通信経路の再設定をさらに含む、請求項16に記載のシステム。

【請求項18】 前記ノード間の動作可能な接続を判断するネットワークモニタ (702) と、稼働中のノードのための情報を処理する高信頼ユーザプロトコルと、前記動作可能な接続に基づいて通信を再設定する論理ネットワークと、をさらに含む、請求項1に記載のシステム。

【請求項19】 前記記憶装置は、各ノードの各ディスク上の情報の一部のみを記憶する、請求項1に記載のシステム。

【請求項20】 各ノードの各ディスクは、他のディスク上の情報の特性を示す情報もまた記憶する、請求項19に記載のシステム。

【請求項21】 冗長分散型サーバであって、

分散型計算ノード (100, 102, 104, 106) のアレイを含み、前記計算ノードの各々は各他ノードとは異なる記憶された情報を有し、前記記憶された情報は、いずれか1つのノード内の記憶された情報が他のいずれか1つのノード内の記憶された情報と組合せられたときに記憶されたデータを再構築することができるよう、前記計算ノード間で冗長となっており、さらに

前記計算ノードのアレイに接続され、前記計算ノードのアレイ間に冗長通信経路を提供し、どのような所定数のネットワーク障害もシステムの残りのノードの動作に影響を及ぼすことのないように動作する、スイッチングシステム (110, 112) を含み、

前記計算ノードの各々はネットワークステータスを判断する同じプロトコルを実行し、よって、各前記計算ノードが同じネットワーク履歴を見るようにし、各前記計算ノードは、ネットワークシステムの動作を妨害し得るシステム機能状態を検出する検出ルーチンを含む、サーバ。

【請求項22】 前記各ノード上に記憶される情報は、所望の情報のすべてではない一部のみを記憶し、どの2つのノードも同じ情報を記憶することはない、請求項21に記載のサーバ。

【請求項23】 前記記憶される情報は、情報部分および冗長部分を含み、前記冗長部分は他のノードのための情報部分のみを示す情報である、請求項22に記載のサーバ。

【請求項24】 前記冗長部分は、複数の前記ノードが前記情報部分を形成するようアレイへと配列されているアレイ符号から形成され、かつ、前記冗長部分は、前記アレイの対角線方向に沿ったチェックサムで形成される、請求項23に記載のサーバ。

【請求項25】 冗長性を提供するようネットワークを操作する方法であって、

集合的にシステムデータを記憶する複数のノード (100, 102, 104, 106) から分散された読出しを実行する、制御プロセスを実行するステップを含み、各ノードは、生のデータおよび、前記各ノード以外のノード内に記憶され

ている生のデータを示す冗長データを記憶して、いずれか1つのノード内の記憶情報が他のノードにおけるデータと組合せられたときに記憶されたデータを再構築することができるようにし、各ノードは、ネットワークシステムの動作を妨害し得るシステム機能状態を検出する検出ルーチン（702）を含み、さらに分散された読出しを行なうステップを含み、該読出しを行なうステップは、ノードの可用性に関するパラメータを判断するステップおよび、前記パラメータが可用性を示す場合には前記複数のノードから前記生のデータを読出し、また、前記パラメータが可用性を示さない場合には前記複数のノードよりも少ない数のノードから前記生のデータおよび前記冗長データを読出すステップを有する、方法。

【請求項26】 誤り訂正符号を示す生の情報および冗長情報を複数の情報ノードに記憶するステップと、

前記情報ノードのユーザビリティを示すパラメータを判断するステップと、

前記パラメータが前記複数のノードが使用可能であることを示す場合には前記複数のノードから前記生の情報を読出し、かつ、前記パラメータが前記複数のノードのうち少なくとも1部分が使用可能ではないことを示す場合には前記複数のノードよりも少ない数のノードから前記生のデータおよび前記冗長データの両方を読出すステップとをさらに含む、請求項25に記載の方法。

【請求項27】 ノードからの情報を表わす、アレイの各列を形成することによって、情報のアレイを形成するステップと、

データを示す生の情報を含む、各列の生の部分を形成するステップと、

冗長性に関する情報を示す、冗長情報を形成するステップとをさらに含み、前記冗長情報は、前記各ノード以外の他のノードに関する情報を示し、該情報は、前記他のノードから情報を得る特定の形状のエンベロープに沿って取られる、請求項25に記載の方法。

【請求項28】 前記エンベロープは、前記アレイの縁を超えて他ノードに延びる対角線である、請求項27に記載の方法。

【請求項29】 アレイの列に各ノードをマッピングするステップと、

前記アレイの列から冗長情報の2つの行を形成し、前記2つの行を前記列内に

位置付けることによって、 $N-2 \times N$ 個の情報シンボルおよび $2 \times N$ 個の冗長情報シンボルを含む、結果として得られる $N \times N$ のアレイを形成するステップとを含み、前記パリティシンボルは以下の式に従って構築される、請求項25に記載の方法。

【数1】

$$C_{n-2,i} = \sum_{K=0}^{n-3} C_{k,(i+k+2)_n}$$

$$C_{n-1,i} = \sum_{K=0}^{n-3} C_{k,(i-k-2)_n}$$

式中、 $i = 0, 1, \dots, n-1$ および $\langle x \rangle_n = X \bmod n$ である。

【請求項30】 該記憶装置はビデオ情報を記憶する、請求項1に記載のネットワーク。

【手続補正2】

【補正対象書類名】 明細書

【補正対象項目名】 0008

【補正方法】 変更

【補正内容】

【0008】

この場合、いずれかの接続またはバスにエラーがあっても、他方のバスを介して正常な動作を続けることができる。2つの冗長バスおよび2つの冗長サーバを有するシステムを、デュアルバス、デュアルサーバと呼ぶ。このようなデュアルバス、デュアルサーバシステムは、単一のネットワーク障害は許容することができる。しかし、このようなシステムでは、通常、各サーバ上にすべての情報の複製を作成しなければならない。

システムの冗長性を提供する試みは、「フレキシブルなプロトコル能力を有する高速スイッチングシステム」と題された公報EP-A-0366935号に記載されているシステムを含む。この文献は、複数のスイッチング面を含んで、1つのスイッチング面またはリンクが障害を起こした場合にも、残りのスイッチング面を使用することによって稼動し続ける、高速のスイッチングシステムを開示

している。このスイッチング面は特定の種類の、特定の目的のためのものである。さらに、このD 1の文献は、システムの再構成およびデータの復元に関するものではない。

【国際調査報告】

INTERNATIONAL SEARCH REPORT

| | |
|--|---|
| International Application No. PCT/US 98/20532 | |
| A. CLASSIFICATION OF SUBJECT MATTER IPC 6 G06F15/173 H04N7/173 G06F11/00 | |
| According to International Patent Classification (IPC) or to both national classification and IPC | |
| B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) IPC 6 G06F H04N | |
| Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched | |
| Electronic data base consulted during the international search (name of data base and, where practical, search terms used) | |
| C. DOCUMENTS CONSIDERED TO BE RELEVANT | |
| Category * | Citation of document, with indication, where appropriate, of the relevant passages Relevant to claim No. |
| A | M.M. BUDDHIKOT ET AL.: "Design of a large scale multimedia storage server" COMPUTER NETWORKS AND ISDN SYSTEMS, vol. 27, no. 3, December 1994, pages 503-517, XP002093312 The Netherlands see page 508, right-hand column, line 1 - line 16 --- -/-- |
| <input checked="" type="checkbox"/> Further documents are listed in the continuation of box C. | |
| <input checked="" type="checkbox"/> Patent family members are listed in annex. | |
| * Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier document but published on or after the international filing date "L" document which may throw doubt on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "Z" document member of the same patent family | |
| Date of the actual completion of the international search 12 February 1999 | Date of mailing of the international search report 26/02/1999 |
| Name and mailing address of the ISA European Patent Office, P.B. 5818 Patentkanal 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040. Tx. 31 651 epo nl Fax (+31-70) 340-2016 | Authorized officer Absalom, R |

Form PCT/ISA/210 (second sheet) (July 1992)

INTERNATIONAL SEARCH REPORT

International Application No.
PCT/US 98/20532

| C. (Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT | | |
|---|--|-----------------------|
| Category * | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
| A | N.J. BODEN: "Myrinet: A Gigabit-per-Second" IEE MICRO, vol. 15, no. 1, February 1995, pages 29-36, XP000501486 los alamos, ca, usa cited in the application see page 31, right-hand column, line 12 - line 16; figure 2 see page 34, right-hand column, line 4 - line 17 | 1 |
| Y | US 5 579 475 A (BLAUM ET AL.) 26 November 1996 cited in the application | 28, 29 |
| A | see the whole document | |
| X | P.C. WONG ET AL.: "Redundant Array of Inexpensive Servers (RAIS) for On-demand Multimedia Services" IEEE INTERNATIONAL CONFERENCE ON COMMUNICATIONS, vol. 2, 8 June 1997, pages 787-792, XP000742048 Montreal, Canada | 27 |
| Y | see the whole document | 28, 29 |
| A | US 5 612 897 A (REGE) 18 March 1997 see the whole document | 1-30 |
| X | US 5 630 007 A (KOBAYASHI ET AL.) 13 May 1997 see the whole document | 25, 26 |
| A | WO 91 14229 A (SF2 CORPORATION) 19 September 1991 see page 1 - page 6, line 24; figure 1 | 1, 30 |
| A | S. NAKAMURA ET AL.: "Distributed RAID style video server" IEICE TRANSACTIONS ON COMMUNICATIONS, vol. e79-b, no. 8, 1908 - August 1996, pages 1030-1038, XP000628640 Japan | |
| A | A. COHEN ET AL.: "Segmented Information Dispersal (SID) for fault-tolerant video servers" PROC. OF THE SPIE, vol. 2604, 23 October 1995, pages 58-69, XP000578588 usa | |

Form PCT/ISA210 (continuation of second sheet) (July 1992)

INTERNATIONAL SEARCH REPORT

International Application No.

PCT/US 98/20532

C. (Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

| Category * | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|------------|---|-----------------------|
| A | R. TENARI: "High Availability in Clustered Multimedia Servers" PROC. 12TH CONFERENCE ON DATA ENGINEERING, 26 February 1996, pages 645-654, XP000632617 New Orleans, usa | |

I

Form PCT/ISA/210 (CONTINUATION of second sheet) (July 1992)

INTERNATIONAL SEARCH REPORT

Information on patent family members

International Application No.

PCT/US 98/20532

| Patent document cited in search report | Publication date | Patent family member(s) | Publication date |
|---|---------------------|----------------------------|---------------------|
| US 5579475 A | 26-11-1996 | US 5271012 A | 14-12-1993 |
| | | DE 69408498 D | 19-03-1998 |
| | | EP 0632376 A | 04-01-1995 |
| | | JP 2750316 B | 13-05-1998 |
| | | JP 7028710 A | 31-01-1995 |
| | | EP 0499365 A | 19-08-1992 |
| | | JP 2514289 B | 10-07-1996 |
| | | JP 4310137 A | 02-11-1992 |
| | | EP 0519669 A | 23-12-1992 |
| | | US 5351246 A | 27-09-1994 |
| US 5612897 A | 18-03-1997 | NONE | |
| US 5630007 A | 13-05-1997 | JP 8329021 A | 13-12-1996 |
| | | GB 2299424 A, B | 02-10-1996 |
| WO 9114229 A | 19-09-1991 | US 5388243 A | 07-02-1995 |
| | | AU 7486091 A | 10-10-1991 |
| | | CA 2077447 A | 10-09-1991 |
| | | EP 0518965 A | 23-12-1992 |

フロントページの続き

| (51)Int.Cl. ⁷ | 識別記号 | F I | ターマコード (参考) |
|--------------------------|--|---------------|-------------|
| H 0 4 L 12/28 | | H 0 4 N 7/173 | 6 2 0 D |
| H 0 4 N 7/173 | 6 2 0 | H 0 4 L 11/20 | C |
| (81)指定国 | EP(AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG), AP(GH, GM, KE, LS, MW, SD, SZ, UG, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, UZ, VN, YU, ZW | | |
| (72)発明者 | ボーシアン, バスケン カナダ、エイチ・4・ジェイ 1・エック ス・4 ケベック州、モントリオール、 11675、ラビン・アパートメント・6 | | |
| (72)発明者 | ファン, チェンゴン アメリカ合衆国、91106 カリフォルニア 州、バサデナ、エス・メリディス・アベ ル、156、アパートメント・115 | | |
| (72)発明者 | レマイエウ, ボール アメリカ合衆国、91106 カリフォルニア 州、バサデナ、イー・デル・マー・ブール バード、1032、アパートメント・301 | | |
| (72)発明者 | リーデル, マーカス・デイビッド・ダニエル カナダ、ジェイ・4・エックス 1・テ ィ・2 ケベック州、ブロッサード、スト ラウス、1015 | | |
| (72)発明者 | シュ, リハオ アメリカ合衆国、91106 カリフォルニア 州、バサデナ、サウス・ウィルソン・アベ ニュー、307、アパートメント・5 | | |